

AN IMPROVED SOURCE MODEL FOR
A LINEAR PREDICTION SPEECH SYNTHESIZER

BY
TUNG-CHU HU

A DISSERTATION PRESENTED TO THE GRADUATE SCHOOL
OF THE UNIVERSITY OF FLORIDA IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

UNIVERSITY OF FLORIDA

1985

To my parents,
and to my wife

ACKNOWLEDGMENTS

I am deeply grateful to my advisor and committee chairman, Dr. Donald G. Childers, who encouraged me to explore my ideas in the field of speech science. Throughout my investigation, Dr. Childers provided guidance, assistance and financial support. I would like also to thank Dr. R. W. Couch II, Dr. F. J. Taylor, Dr. J. C. Principe and Dr. H. C. K. Yang, for serving on my supervisory committee and advising me on various aspects of this dissertation. My colleagues under Mind-Machine Interaction Research Center helped me in many ways. I appreciate their interest and contributions to my research. Finally, I would like to dedicate this dissertation to my wife, Yuh-Yuh, who patiently shared every struggle I experienced during the process, and to my parents, who never hesitated to deliver their love, encouragement and support.

TABLE OF CONTENTS

	PAGE
ACKNOWLEDGMENTS	iii
ABSTRACT	iv
CHAPTERS	
1 INTRODUCTION	1
1.1 Speech Production Mechanism	1
1.1.1 Respiratory System	2
1.1.2 Acoustic Mechanism	3
1.2 Previous Research on Speech Production	3
1.3 Models for Speech Synthesis	4
1.3.1 Formant Model	5
1.3.2 Source-Filter Model	6
1.3.2.1 Articulatory synthesizer	7
1.3.2.2 Formant synthesizer	7
1.3.2.3 LP synthesizer	10
1.3.2.4 Comments on the three types of synthesizers	14
1.4 Research Issues and Objectives	14
1.5 Description of Chapters	17
2 SOURCE PROPERTIES	20
2.1 Review of Existing Acoustic Measures	21
2.1.1 Perturbation Measures	21
2.1.2 Characteristics of the Global Flow Waveform	21
2.1.2.1 Quantitative analysis based on parameters of source models	22
2.1.2.2 Spectral tilt	22
2.1.3 Vocal Noise	22
2.1.4 Basis of the Inverse Vocal Tract Filter	23
2.1.5 Vocal Intensity	23
2.1.6 Remark	24
2.2 Global Inverse Filtering	24

1.3 Continuous-frequency Analysis and Differentiated Glottal Flow	38
1.4 Choice of Model Type	39
1.5 Data Collection and Methodological Considerations	39
1.5.1 Experimental Data Base	39
1.5.2 Vocal Quality	39
1.5.3 Analytical Layers	39
1.5.4 Standardization of Pitch Period	40
1.6 Feature Extraction	43
1.6.1 Formant Measures	44
1.6.2 Spectral Tilt	45
1.6.3 Glottal Phase Characteristics	45
1.6.3.1 Glottal properties	45
1.6.3.2 Absorption index	46
1.6.4 Vocal Noise	46
1.6.4.1 Noise estimation	46
1.6.4.2 Properties of vocal noise	46
1.6.4.3 Brief summary	47
1.7 Discussion	47
1.8 Conclusion	49
2 SOURCE MODELING	71
2.1 Review of Previous Research	71
2.2 Excitation Source	74
2.2.1 Vowel Segments: Excitation Pulse	74
2.2.1.1 Vector quantization	75
2.2.1.2 Maximum-entropy algorithm	81
2.2.1.3 Cluster splitting	81
2.2.1.4 Codebook training	83
2.2.2 Unvoiced/voiceless Segments: White Noise	85
SPEECH ANALYSIS/SYNTHESIS EVALUATION	89
4.1 Analysis Scheme	90
4.1.1 Orthogonal Covariance Method	91
4.1.2 VQPS Classification	95
4.1.3 Identification of Glottal Closure Interval (GCI)	95
4.1.4 Codeword Sampling	99
4.1.4.1 Vowel extraction: glottal codebook	99
4.1.4.2 Unvoiced-extraction: stochastic codebook	101
4.2 Synthesis Scheme	103
4.2.1 Interpolation of Glottal Phase	105
4.2.2 Interpolation of LP Coefficients	110
4.2.3 Spectral Flattening	113
4.2.4 Effect of Vocal Noise	113
4.2.5 Source-train Iteration	113

4.3.6 Generation of Glottal Impulse	127
4.3 Class Determination	127
4.3.1 Case of Nasal Excitation, A_0	127
4.3.2 Case of Unvoiced Excitation, A_0	128
4.3.3 Voicing Transition	129
4.4 Subjective Quality Evaluation	129
CONCLUDING REMARKS	130
5.1 Summary	130
5.2 Possible Improvements	132
5.2.1 Estimation of Vocal Noise	132
5.2.2 OCF Identification	133
5.2.3 Excitation Source	133
5.2.4 Ripple Effect	134
5.2.5 Sampling Resolution	135
5.2.6 Spectral Estimation	136
5.3 Applications	136
5.3.1 Quality Monitor	136
5.3.2 Speech Coding	136
5.3.3 Voice Conversion	136
5.3.4 Text-to-Speech Synthesis	137
REFERENCES	138
BIOGRAPHICAL SKETCH	140

Abstract of Dissertation Presented to the Graduate School
of the University of Florida in Partial Fulfillment of the Requirements for the
Degree of Doctor of Philosophy

**AN IMPROVED SOURCE MODEL FOR
A LINEAR PREDICTIVE SPEECH SYNTHESIZER**

By

Thao/Thu Ha

May, 1990

Chairman: Dr. D-G. Chaffin
Major Department: Electrical Engineering

Though some progress has been made towards producing natural-sounding speech using the linear predictor (LP) technique, an appropriate feature-based parametric representation has not been fully developed for this type of synthesizer. The intent of this research is to verify the importance of selected acoustic measures by means of LP analysis within the source-filter theory, and then use the derived information to develop an LP synthesizer that is capable of synthesizing high-quality natural-sounding speech.

In order to carry out the requirements of this research, we divide the relevant issues into two separate but related phases. In the first phase we propose methods for isolating and extracting the acoustic features of vocal quality. Based upon a comprehensive speech production model, the LP analysis is used to estimate spectral properties of the speech signals and of the glottal source. Relevant source features, such as time, pitch and duration, are recovered from the synthesized signals. An algorithm is developed to extract the time domain characteristics of the vocal noise. Various aspects of each extracted noise are examined subsequently. To illustrate the above mentioned analysis techniques, the measured acoustic parameters of the proposed speech model for three voice types (male, female and

locally) are provided as representative examples. It is anticipated that our findings will contribute to the understanding of the problems of modeling the excitation source and the LP system.

In the second phase we propose a novel source model to generate the excited signal in terms of the glottal phase characteristics. Depending on the voicing condition of the analyzed speech, the excitation source is formulated as two separate codebooks, i.e., a glottal codebook for voiced segments and a stochastic codebook for unvoiced and silence segments. Methods for determining voicing intervals are presented, along with procedures for searching the codebooks for the appropriate excitations. Short pitch spectrograms schemes are proposed for speech synthesis, we describe procedures for identifying the instants of glottal closure and for interpolating the excitation pulses as well as the LP coefficients. Moreover, we account for the effects of nasal noise and source-filter interaction, which are generally ignored in most synthesizers. Finally, a method for determining the missing gain is given. This method also serves as a supplementary tool to evaluate the relationship between the gain and power spectrum. Informal listening tests were used to evaluate the speech processing techniques. The listening tests revealed that the quality of synthetic speech was close to that of the original speech. The results indicate that our source model is able to characterize the glottal features and that the overall speech production model is quite adequate for high-quality synthesis.

CHAPTER 1 INTRODUCTION

Speech is a sophisticated skill that humans have developed for efficient communication. This skill involves not only linguistic information but also acoustic factors that convey the speaker's identity and other aspects of the speaker's physical and emotional state. Although our current knowledge is insufficient to reveal the biological codes that describe the phonatory system, the mechanism of speech production is becoming comprehensible due to the advances in acoustic theory and computing technologies. Phonatory acoustics forms the basis for all present-day speech synthesizers. The increasing use of speech synthesizers in the marketplace has produced great demand for products that can generate "high-quality" speech. In fact, quality degradation with regard to existing speech synthesizers mostly results from unwanted-sounding characteristics, which are known to cause perceptual difficulties (Peters and Houtgast, 1982; Paoli et al., 1983). The aim of this dissertation was to seek methods that could improve the naturalness of synthetic speech. We restricted our attention to the Linear Prediction (LP) technique, because of its simplicity and accuracy in speech processing.

Following a brief introduction to speech production, we give an overview of some existing synthesis techniques. Then we provide the details of speech modeling and processing. This overview, associated with the basic knowledge of the phonatory system, describes our explanation of some efforts made to improve the naturalness of synthetic speech.

1.1 Speech Production Mechanism

Perhaps the easiest way to describe the speech production mechanism is to explain the physiological function of the anatomy of the human vocal system. In general, the speech

production system can be divided into two systems, namely, the excitation source and acoustic modulation, which shape the modulated spectrum to form intelligible sounds (or phonemes).

3.1.1 Excitation Source

The lungs act as an air reservoir, expelling air up the trachea to the vocal folds. During periods of voiced speech, the vocal folds open and close in a quasi-periodic fashion, producing a pulsating stream. During the periods of unvoiced speech, the vocal folds are held apart so that the airstream is less disturbed and can be considered as a nearly turbulent source. The vocal folds have an important role in determining the characteristics of vocal quality. The modulation of the vocal folds can be described by the aerodynamic-mechanical theory (Jiang et al., 1997; Jiang, 1998). Basically, the motion of the vocal folds is controlled by several interplaying forces that cause the abduction and adduction of the folds. When the subglottal pressure builds up to a certain level, the vocal folds are pushed apart and the air is then released through the glottis. The volume velocity of air passing through the glottis increases as the vocal folds keep opening. As the velocity increases beyond some threshold, pressure across the folds begins to drop and then results in a Bernoulli effect. This effect, in conjunction with the elastic resistance of the folds, induces the adduction of the vocal folds at the time these two effects outweigh the subglottal pressure. When the vocal folds close, the subglottal pressure builds up again and the entire procedure repeats. Such a repetitive cycle is referred to as a glottal period/its reciprocal is denoted as the fundamental frequency.

Noise generated by turbulence is another important source of speech production. The airflow emerging from the lungs can cause turbulent streaming while passing through a vocal aperture, which is either the vibrating vocal folds or a constriction along the vocal tract. Such turbulence causes either vocal apertures open sufficiently or the airflow decreases. The possibility of turbulent flow is reflected by the value of the Reynolds number, which

characterizes the viscosity of the airstream as either laminar turbulent or somewhere in between (Adams, 1980). With these kinds of characteristics, there is no doubt that the substance becomes an essential element in fricatives, aspirates, glides, whistles and finally vowels. This fact necessitates the use of a vowel source while synthesizing these particular sounds.

1.1.1 Acoustic Modelman

The human vocal tract, extending from the glottis to the lips, can be considered as an acoustic tube of nonuniform shape varying as a function of time. Components that lead to this time-varying change include the lips, jaw, tongue, velum and nasal cavity. During the periods of nasalized vowels, the velum closes-off the nasal tract from the vocal tract. Thus, the acoustic tube only exhibits poles in its transfer function. When the velum is lowered, the vocal tract is acoustically coupled with the nasal tract, forming a pole-zero system. As the tube varies the shape for different sounds, the resulting transfer function is such that it emphasizes certain frequency components of the glottal wave and/or de-emphasizes others. The resonant peaks of the speech output due to the poles are referred to as formants, whereas the valleys due to the zeros are referred to as anti-formants.

1.1.1 Previous Research on Speech Production

The earliest efforts of speech research were directed to exploring the physiological nature of the human phonatory system. At first time, the speech synthesizers played a fundamental role in learning the process of speech production. The talking machine, designed by van Kampelen in 1793, consisted a bellows which supplied air to a reed (Shangpa, 1978), the bellows and the reed were obviously used to simulate the lungs and the vocal folds respectively. A hand-cranked mechanism was provided to simulate the acoustic response of the vocal tract. This machine was reported to produce only a few vowels. Modern speech synthesizers are directed to mimic the technologies developed over this

century have come out with sophisticated techniques which greatly improved the quality of synthetic speech. Such a technological evolution accompanied with the emerging understanding of speech acoustics, gradually shifted the focuses and interests of speech synthesis to other applications. The most significant influence was the Voronoi awarded by Dudley (Dudley, 1939), whose efforts spawned a sub-field of communication engineering. Research in this sub-field was aimed at the efficient encoding and transmission of speech information. The techniques of interest were directed toward obtaining acceptable quality at low bit rates, using reasonable computational resources in a real time environment. Research issues encompassed methods to improve quality, robustness, delay and complexity.

As speech synthesis techniques have continued to improve in recent years, many speech synthesizers have been employed to implement voice response systems for computers, which are called the "text-to-speech" techniques. Speech synthesis is the issue of "text-to-speech" means automatically producing voice response according to a text input. The capability of voice response offers possibilities for automatic information services, computer based classrooms, talking aids for the visually handicapped, and reading aids for the visually impaired.

1.1 Models for Speech Synthesis

This research is directed toward improving the speech production model. We define the term "speech analysis" as the procedure used to extract the speech production model parameters from the speech signal and "speech synthesis" as the procedure used to reproduce the synthetic speech signal by controlling and updating the appropriate parameters obtained from the speech analysis.

Modern speech synthesizers can be classified into two groups, one based on the Fourier transform methods and the other based on the linear source filter model.

3.3.1 Fourier Model

The Fourier transform has traditionally been used to study speech signals because it provides a frequency-domain analysis of the phonatory and auditory properties of speech signals. Using the Fourier model, the speech signal is analyzed using short-time Fourier analysis (STFA), while synthesis is carried out by an inverse transform (Allen, 1977, Allen and Rabiner, 1977). The term “short-time” implies that the speech spectrum is stationary over a short interval of time. This is a valid approach to speech processing because many psychoacoustic and physiological studies have shown that the human ear performs a type of short-time spectral analysis of acoustic signals.

The channel vocoder is the oldest form of speech coding device that employs Fourier analysis and synthesis (Dudley, 1998). This vocoder is constituted by several bandpass filters, each of which is employed to preserve the magnitude Fourier transform of the speech signal within a specific band. An additional channel is needed to transmit other information regarding the excitation, e.g., the voiceless/voiced signal and the pitch period for voiced speech. Consequently, the concept of source excitation was incorporated into the configuration of the channel vocoder.

Another Fourier-based model that has experienced popularity is the phase vocoder (Pnagen and Golden, 1960). The major success of this technique originates from a polar representation of the Fourier transformation, i.e., phase and amplitude, which leads to an overview of information bandwidth. Unlike the channel vocoder, which neglects the phase spectrum, the phase vocoder exploits the phase information through the derivative of the phase spectrum. Furthermore, it provides flexibility for expanding and compressing the time scale through the manipulation of the instantaneous frequency. Emerging from a similar idea, a new class of models called “alluvial vocoders” were developed and have proliferated since the early 1980s (Riedler, 1981; Klapuri and Erkkila, 1982; Almeida and Silva, 1984; McAdams and Quatieri, 1984; Tzaneto et al., 1990). For such vocoders, the speech signal

within each frame is approximated by a superposition of sinusoids with time-varying amplitude and frequency:

$$s(t) = \sum_{i=1}^N a_i(t) \cos(\omega_i t + \phi_i(t)), \quad 0 \leq t < T \quad (2-1)$$

where N is the number of sinusoids, $a_i(t)$ is the amplitude of the i th sinusoid, $\phi_i(t)$ is the corresponding frequency and T is the frame length. The variation of amplitudes, a_i 's, and phases, ϕ_i 's, within a short interval is usually described by first- and third-order polynomials respectively as

$$a_i(t) = A_i + (t/T)B_i \quad (2-2)$$

$$\phi_i(t) = \omega_i t^2 + \tau_i t^3 + \tau_{i2} t + \tau_{i3} \quad (2-3)$$

These polynomials are then applied to an interpolation rule for the instantaneous values of amplitude and phase as well as frequency. With a little abuse, speech quality obtained using this model is virtually indistinguishable from the original.

Among the sinusoidal coders, approaches for processing the Fourier-based parameters can be divided into two classes. Members in one class separate the pitch harmonics from the spectrum envelope, and only apply the sinusoidal processing techniques to the harmonics. In other words, the a_i 's in Eq. (2-1) are obtained by other means such as linear prediction and cepstrum analysis. Members in the other class, on the other hand, consider the a_i 's as part of the results of Fourier analysis. Interestingly the generation of noise excitation also exhibits two different forms, i.e., either white noise or a signal with random phases but constant amplitudes.

1.3.3 Source-Filter Model

The source-filter model was developed by Fant in the late 1950s (Fant, 1975; Fant, 1960). In this model the speech signal is modeled as the filtered output of an excitation source

by quasi-periodic pulses for voiced speech or by random noise for unvoiced speech. The transfer function of the network is defined as the ratio of the Laplace transform of the sound pressure from the lips of the speaker to the volume velocity of the airflow passing the vocal folds. In the case of speech production, a speech signal indicates both characteristics of the source and the network. Formed based on the source-filter theory, speech synthesizers can be further classified into three categories, namely, the L-F formant and articulatory synthesizers.

1.1.2.1 Articulatory synthesizer

The articulatory synthesizer is a direct approach that simulates speech production and propagation from the viewpoint of anatomy and physiology. In order to describe the wave propagation by means of aerodynamic equations, we must specify such parameters as subglottal pressure, density of the vocal folds and viscosity of the vocal tract, in addition to the movement of the articulator coordinates and the changes in the vocal fold configuration. As shown in Figure 1-1, although the overall computation can be broken down into a sequence of subproblems of constant cross-sectional areas, the complexity involved in this aerodynamic and mechanical system is still considerable. Therefore, researchers have attempted to convert the gross features of vocal fold vibrations into a model with acoustic parameters. Likewise, the vocal and nasal tracts were represented by an equivalent circuit such as an acoustical transmission line (Figure 1-2). Furthermore, the control system for driving the area functions of the vocal tract was developed by matching the formant characteristics of the model to those of real speech.

1.1.2.2 Formant synthesizer

The development of the formant synthesizer is mainly based on the perceptual characteristics of the human auditory apparatus. In this type of synthesizer, the transfer

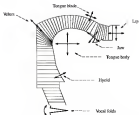


Figure 1-1. Articulatory model of human vocal tract and the associated control variables.

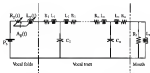


Figure 1-2. Equivalent circuit for the vocal system.

function is directly controlled by the use of resonant and anti-resonant filters whose center frequencies and bandwidths can be individually specified. The resonant filters can be connected either in parallel or in series so as to distribute the production of both nasal and nonnasal sounds. An excitation model mimicking the natural excitation is used to provide the source properties. Figure 1-3 shows the design of a typical formant synthesizer.

1.3.3.3 LP synthesis

The LP synthesizer consists of an excitation source and a time-varying all-pole filter (Figure 1-4). The all-pole filter determines the spectral envelope of the synthesized speech, and the excitation source provides the fine structure of the spectrum harmonics. The all-pole filter is derived from a mathematical approach that regards the speech signal as an autoregressive process, that is, the current sample is a linearly weighted sum of previous samples. This approach yields an accurate and efficient representation of the short-time spectrum of speech signals. Since the human ear is mostly sensitive to the magnitude spectrum of an acoustic signal, the ability of preserving the spectral envelope by the LP analysis is the main reason for its success.

The representation for the spectral envelope based on LP analysis also has many applications to other types of vocoders. For instance, the cepstrum, which is obtained from a homomorphic system (Figure 1-1), is an alternative form that modifies the short-time speech spectrum (Oppenheim, 1983). The impulse response, h_n , computed from the cepstrum, can be considered as the coefficient sequence of an FIR filter exhibiting a similar spectral envelope of the all-pole filter:

$$H(z) = \sum_{n=0}^N h_n z^{-n} = \frac{G}{1 - \sum_{k=1}^K a_k z^{-k}} \quad (1-6)$$

where a_k is the LP coefficient, and G is the gain. The frequency spectrum applied to the source is modified by a correlation between the sequence, h_n , and the excitation, e_n .

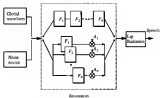


Figure 1-3. Block diagram of formant synthesis.

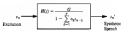


Figure 1-4. Block diagram of LP synthesizer



Figure 1-5 Block diagram of homomorphic system (a) analysis process
(b) synthesis process

1.1.1.4 Comments on the three types of synthesizers

The advantages and disadvantages of the three types of synthesizers are given in Table 1-1. For semi-LP synthesizers, the deficiency of the all-pole filter is ameliorated by employing a pole-zero model (Liu and Schroeder, 1978; Childers et al., 1981). Also, the excitation function may be replaced by synthesized pulses or innovative resonances that contain the desired signal. (This is discussed in detail in Section 3.1.1.) Moreover, an independent control of the spectral characteristics is achieved by factoring the filter into resonators and anti-resonators (Karnikian, 1984; Childers et al., 1989c). For some formant synthesizers, the dynamics of spectral characteristics are enhanced simply by inserting formants into the glottal source (Pattis and Ljungqvist, 1984). Likewise, the effect of vowel-vowel transitions may be simulated by either modifying the glottal wave shape or the formant bandwidths or by incorporating a control signal (Guthrie et al., 1976; Yin et al., 1980; Pattis and Ljungqvist, 1984; Wang, 1991).

Since each of the above-mentioned schemes increases the computational burden of the processing unit, the complexity is no longer a major drawback only for the articulatory synthesizer. Besides, many articulatory synthesizers/formants of the articulators are determined by comparing the formants of the synthetic speech with that of the original speech (Pattis and Ljungqvist, 1984; Prato et al., 1982). Consequently, one type of source-filter synthesizer is not particularly different from the others.

1.4 Research Issues and Objectives

Ideally, a speech synthesizer should have the ability to produce any desired voice quality. From this standpoint, the criterion of voice quality are directly related to the control parameters of a synthesizer. In other words, these control parameters may modify the quality attributes that the human ear uses to discriminate voice types. This use of a speech production model for speech research is called "analysis-by-synthesis."

Table 1-1. Components for articulatory, formant, and LP synthesizers.

	LP synthesizer	Formant synthesizer	Articulatory synthesizer
Advantage	<ol style="list-style-type: none"> 1. few parameters are required. 2. fast algorithms are available. 3. synthetic speech is intelligible at a rate as low as 2 B/s. 	<ol style="list-style-type: none"> 1. glottal waveform and formant frequencies can be controlled independently 2. control parameters correlate with the acoustic aspects of the speech sounds. 3. source-filter interaction can be simulated by modifying the glottal source. 	<ol style="list-style-type: none"> 1. control parameters are directly related to the articulatory mechanisms. 2. source-filter interaction can be modeled. 3. articulatory parameters can be manipulated prior to signal conversion to the speech signal.
Disadvantages	<ol style="list-style-type: none"> 1. nasals, fricatives and many consonants cannot be properly produced. 2. low pitch values often sound hoarse. 3. control parameters show little relation to the anatomy and physiology of speech production. 4. source-filter interaction cannot be produced in a direct manner. 	<ol style="list-style-type: none"> 1. formant-based formant conversion is difficult. 2. synthetic speech may sound too smooth. 	<ol style="list-style-type: none"> 1. explaining exactly the processes of interactions and voluntary patterns of vocal folds are difficult. 2. considerable computation is required.

Our ultimate objective in this dissertation was to develop a high-quality speech synthesizer. The quality of speech, in general, is referred to as the total auditory impression the listener experiences upon hearing the speech of a speaker. It consists of two factors, namely, naturalness and intelligibility. Through the progressive understanding of speech production, researchers should be able to validate the hypothesis that the intelligibility of speech signals depends largely on the vocal tract, while the source characteristics determine the naturalness of the voice. The intelligibility aspect concerns how since most people (lay speakers) are capable of conveying the intended speech content correctly. Instead, we are more interested in the vocal source because of its contribution to the naturalness of speech. For this reason, the words "quality" and "naturalness" will be considered equivalent in this dissertation.

To accomplish our objective, we decided to divide the research issues into two separate but related phases. In the first phase we discussed how to obtain accurate measures by LP techniques. Three types of voiced speech (nasal, vocal fry, breathy) were used as representative examples to illustrate the source properties.

We have selected the LP technique to accomplish this study despite some arguments against the use of such a technique. The LP analysis, in our opinion, is more than adequate because source properties are all extractable from the residual signal obtained by inverse filtering of the speech signal. This argument becomes clear in Chapter 3 when we discuss the relationship between the residual and the volume-velocity flow. We identify the significance of acoustic measures extracted from both the excitation and speech signals and subsequently correlate these measures to the control parameters of a speech production model.

The knowledge gained in the first phase of the research is useful for the design of a source model for an LP synthesizer. It has long been known that the lack of glottal characteristics is the primary reason leading to the poor quality of LP synthesizers. In the second phase of the research, we first try to develop a source model to simulate the residual

signal. Such a source model will be presented in the form of a cookbook that will be incorporated into a newly designed speech production model. In addition to source modeling, other factors, such as the incorporation of LP coefficients, turbulent noise, source-mut interaction, etc., have to be taken into consideration. Thus, we will present our methods and strategies to deal with these factors.

The efficacy of the source and speech production models is determined by evaluating the quality of synthetic speech. We have taken the analysis-by-synthesis approach in studying speech quality. This approach provides information about whether the important acoustic features are successfully maintained during the modeling process. While an objective quantitative measure is available for performing speech evaluation, informal subjective listening tests were conducted to assess the quality of the synthetic speech samples. The overall research plan is presented as a schematic diagram in Figure 1-4.

1.3 Description of Chapters

Chapter 2 describes the procedures for measuring vocal source properties by linear predictive analysis. Following a retrospect of some existing acoustic measures, our first focus is on solving the relationships between these measures and the control parameters of a comprehensive speech production model. In particular, under the guidance of this model, we propose methods for identifying and isolating the acoustic characteristics of vocal quality. Three voice types are provided as representative examples to illustrate the proposed model and analysis techniques. Knowledge gained in this chapter contributes to the understanding of general problems of source modeling and speech processing, which we present in Chapters 3 and 4.

Chapter 3 deals with the modeling of the excitation source. Depending on the voicing condition, we divide the excitation into two categories, i.e., voiced and unvoiced. A novel glottal source model is proposed to describe the voiced excitation in terms of the glottal phase transformations, while the nonlinear sequences are used to simulate the unvoiced excitation.

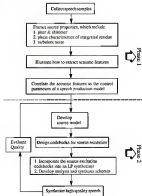


Figure 1-6: Schematic diagram of speech plan.

both types of variations are formulated into codebooks. Our methods of generating the codebooks are motivated with these two types of variations individually.

The linear predictive analysis and synthesis schemes used in this study constituted the first two parts of Chapter 5. Issues such as the voicing decision, Glottal Closure Instant (GCI) identification, vowel-onset searching, vocal noise, source-filter interaction and gain determination are addressed. The overall performance of these schemes is dependent on how closely the reproduced speech resembles the original. While no reliable objective quality measure is currently available, we evaluate the synthesised speech by informal listening tests.

Chapter 5, the last chapter, summarises the results of this study, discusses possible improvements to the proposed model and finally recommends some potential applications.

CHAPTER 2 SOURCE PROPERTIES

A better understanding of speech production is important for the measurement of speech quality as well as for the development of a natural-sounding speech synthesis model. In this chapter we are particularly interested in the global source properties that affect the perceptual quality of the voice. The closeness of the relationship between the excitation source and the residual speech quality requires source-related parameters to describe acoustic and perceptual features, as well as methods to extract the parameters. The analysis-by-synthesis technique is a general approach to speech analysis (Rabiner and Scholer, 1979; Paris, 1983). In principle, we establish the speech production model and then derive the model parameters used to reproduce speech signals. Speech synthesis, in conjunction with perceptual evaluation, plays a role in validating the significance of the acoustic features in terms of the model parameters. As the speech production models become more and more sophisticated, many detailed acoustic/perceptual correlations will be easily verified by the analysis-by-synthesis approach.

Our major concern is focused on efforts that will establish a relationship between model parameters and acoustic features measured from the speech signal. Following a brief review of existing acoustic measures and a background description of source filtering techniques, we discuss the relationship between two commonly reconstructed source excitation signals, namely, the residual signal and the differentiated glottal flow waveform. In order to facilitate the acquisition of the source excitation, we have used the LP technique as a vehicle to complete this research. A new LP synthesis model with appropriate source features is also proposed. Near analogies of these types of phenomena, i.e., model, vocal

ity and locality, were used as representative examples to validate the competence of the proposed model.

2.1 Review of Existing Acoustic Measures

Basically, researchers have used five types of acoustic measures to study vocal quality:

- (1) Perturbation measures,
- (2) Characteristics of the glottal flow waveform,
- (3) Vocal noise,
- (4) Slopes of the longest vocal tract filter,
- (5) Vocal intensity

2.1.1 Perturbation Measures

Voiced speech is generated by the vibration of vocal folds. Abnormal vibratory patterns of vocal folds has long been known to result in abnormal or deviant voices (Moore, 1976). Statistical properties of the cycle-to-cycle variations in voiced speech have proven useful to characterize vocal quality (Ashcraft and Hammarberg, 1986; Schoenberg, 1988; Fitch and Stein, 1988; Riekman et al., 1990). The perturbation ratio (fundamental frequency and amplitude of sustained utterances, termed *pitch* (Lukkarinen, 1987) and *intensity* (Kobayashi, 1988), respectively, were two of the first acoustic measures reported to be correlated with vocal pathology. Since then, other perturbation measures have also been shown to be capable of discriminating pathological from normal voices.

2.1.1.1 Characteristics of the Glottal Flow Waveform

The characteristics of the glottal flow consist either the assessment of speech quality can be further classified into two categories. (1) qualitative analysis based on parameters of source models, and (2) spectral (B).

2.1.2.1 Quantitative analysis based on parameters of source models

Monitoring the glottal flow waveform is a direct means of studying the variations of the glottal source (Hoffman and Wosberg, 1981; Jordan et al., 1982; Poon, 1989). In order to assess such variations on a quantitative basis, a parametric model must often be introduced. One such model that has been widely adopted for quality assessment in recent years is the LP model (Frost et al., 1983; Fujisaki and Ljungqvist, 1989; Frost and Lee, 1991; Gohl, 1993 & 1999; Karlsson, 1993; Aho, 1995; Turpeinen and Reinart, 1999; Chiklen and Lee, 1999). This model is useful because it requires an overall flow consistently encountered differential glottal pulses with a minimum number of parameters, and it is flexible in the extent to which it can match various phenomena.

2.1.2.2 Spectral tilt

In addition to the parametric variations of the source models, the spectral tilt of the glottal flow appears to be characteristic of different voice types (Hoffman, 1984; Holo et al., 1986; Märmann and Engelen, 1993). In fact, the magnitude of the spectral tilt is caused by the rigidity of the closing phase and by the abruptness of the glottal closure. The perceived quality of speech is related to the spectral tilt (Chiklen and Lee, 1999). A steeply declining spectral tilt results in a hot quality, whereas a gradually declining tilt produces a more quality. To achieve a quantitative measure of this aspect of vocal quality, the spectrum of the glottal flow is usually approximated by a three-pole model, or equivalently a two-pole model for the differentiated glottal flow. The coefficients of the three- or two-pole models are then used to indicate the spectral tilt.

2.1.2.3 Vocal Noise

Turbulence at the level of the glottis also contributes vocal quality such as hoarseness and breathiness, which is a proximal expression of laryngeal pathologies (Kiani, 1987; Kiani and Kiani, 1999; Chiklen and Lee, 1999). Methods for measuring the turbulent noise current

of the relative intensity (Kusaka et al., 1984; Palacios et al., 1989), the spectral ratio level and the harmonic-to-noise ratio (Kusaka, 1981; Yamato et al., 1983; Yamato et al., 1984; Kusaka et al., 1984a,b; Mats et al., 1987; Children and Lee, 1991). In most cases, these noise measures were influenced by the spectral content of the analyzed speech. Consequently, better methods are needed so that the global flow waveform can be analyzed more precisely.

2.1.4 Form of the laryngeal Vocal Tone Filter

Another aspect of speech spectra that affects the detection of laryngeal dysfunction has been discussed by Deller and Anderson (1986), who represented the speech signal by the roots of the inverse filter and then applied pattern recognition techniques to discriminate the subjects as either normal or pathological. It was found that the discrimination factors employed in detecting laryngeal behavior was more sensitive to the poles attributable to the glottal source than to the formant structure (Deller, 1987). This technique was later applied to the F0G signal by Smith and Chikara (1993) as a method for detecting laryngeal pathology. They concluded that the LP features of F0G signals were more sensitive to pathology detection than similar features measured from speech signals. Recently, the same task was recast on the pattern analysis of LP coefficients by vector quantization (Chikara and Sun, 1995). Inferences based upon their results were consistent with previous research.

2.1.5 Vocal Intensity

Vocal intensity is less specific in quality assessment (Chikara, 1995; Hoshino, 1994). It is largely sensitive to the perceived quality change for loudness. Since the selected speech samples we used were approximately at the same power level after digitization, vocal intensity was not considered an important factor in our research.

3.1.2. Shortcuts

A rudimentary statistical analysis of acoustic parameters may result in a quality predictor that matches well with the objective evaluation (Jilka et al., 1996; Wills and Swales, 1992; Ekman et al., 1990; Pass and Tiesi, 1994). In addition, measures of higher order may also provide more degrees of freedom in statistical analysis. Using these measures in quality assessment causes difficulties in justifying the significance of each individual measure and their correlations. Pass and Tiesi (1990) made an attempt to unify existing pitch, duration and noise measures, however, no effort was made to sort out the relation between the acoustic measures and the overall parameters of a specific speech production model. This motivated us to explore these relations.

3.2. Global Source Filtering

Global source filtering is a popular and efficient means for manipulating the activities of the global source. It is based on the assumption that the source parameters and the supraglottal loading are separable and that the source properties of the speech production model can be uniquely determined. The principle of source filtering is to obtain the global GSF by eliminating the effects of vocal tract transfer function and lip radiation from the speech signal. Figure 3-1 presents the conceptual source filtering model. Notice that under assumption the sequence of the vocal tract transfer function and lip radiation are reversed because the speech production is assumed to be a linear model.

Common methods for global source filtering center on LP analysis (Gerstoft, 1976; Wang et al., 1979; Matsushita and Riedler, 1949; Childers and Laro, 1944; Kishimoto et al. and Childers, 1958; Mitinsky et al., 1946; Childers and Lee, 1991). Among the various methods, the closed-phase covariance analysis is considered the most reliable because no source-tract interaction is involved. However, the disadvantages of this method are: [I] it needs to locate the closed phase very accurately, and [II] it is only feasible when the closed phase is long enough to accommodate the analysis window.

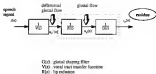


Figure 2-1 Block diagram of global inverse filtering.

Recently, in order to overcome these disadvantages, adaptive approaches have been used to track the rapid change of the parameters of the vocal tract during the glottal phase (Ting and Childers, 1982). In fact, it is more convenient to estimate the composite effect of the glottal pulse, lip radiation and vocal tract together. The vocal tract transfer function could be obtained by removing the source-related roots from the LP polynomial (Childers and Lee, 1983). However, this approach may introduce distortion in the estimated elimination or merging of such roots. Furthermore, since the estimate of the vocal tract parameters is function of each pitch period, the effects of different damping factors caused by the open and closed glottal intervals during a pitch period affect the estimate. Consequently, the estimated glottal flow waveform becomes an "average" waveform for the entire pitch period. This average waveform may not be truly representative of the actual waveform.

From the discussion above, we know that the glottal flow waveform is not always obtainable using the glottal source filtering techniques. However, the estimation of vocal signal is indeed affected by the preceding factors. Moreover, as will be seen in the next section, the removal of the glottal phase characteristics can be achieved from the residual signal. For these two reasons, we concentrate our study on the residual signal. The potential of the residual can be seen from its appearance. It has been observed that the residual extracted from normal vocal consists of periodic sharp spikes and low-level noise components, whereas the residual extracted from deviant vocal exhibits a less distinctive pattern of periodic spikes (Figure 3-5). Because such an observation is not as noticeable as in the speech signal, many researchers advanced the use of the residual signal over the speech signal for the analysis of abnormal voices (Kusler and Makiel, 1975; Sorenson and Horv, 1984; Probst et al., 1987). Recently, the quantitative measures deduced from the residual signal failed to support their claims (Johannsson, 1987). We believe this contradiction is due to the inadequacy of the acoustic measures and the analysis methods. It was noted that the LP coefficients calculated by a fixed frame synchronization method, which was used by

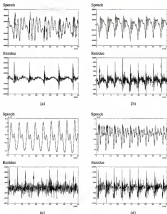


Figure 2-2 Speech and residual waveforms for two normal subjects (a) and (c), and for two pathological subjects (b) and (d). The pathological symptom is hoarseness for subject (a) and a bilateral paralysis of TVC for subject (d).

Schroeder (1962), were affected by the delay and position of the analysed frame. Any small deviation of the estimated coefficients could result in a gross change of the residual signal (Ammari, Peddaz, and Yegorov, 1979). Consequently, the acoustic estimates derived from the fixed-frame autocorrelation method are prone to error. To avoid this problem, a pitch-synchronous-correlation analysis method has been used (Chang and Lu, 1974).

2.3 Correlation between Residual and Differentiated Glottal Flow

It is constructive for us to clarify the relation between the residual and the glottal flow before we explain the characteristics of the glottal source. As shown in Figure 2-1, the inverse filtering can be regarded as a process of uncorrelating the speech signal so as to obtain the excitation waveform. One of the intermediate products is the differentiated glottal flow, while for our purpose the residual signal is the ultimate result. Thus, the correspondence between the residual and glottal flow can easily be illustrated as a filtering process. Here we adopt a two-pole filter to model the spectrum of the differentiated glottal flow. The filter coefficients are obtained by an LP analysis of the modeled LF waveform.

Since the LF model has been successfully used to describe the characteristics of a differentiated glottal flow, we adopt it as an explanatory model for the subsequent discussion. The equations of the LF model are given as

$$E(t) = \begin{cases} E_0 e^{-\alpha_0 t} \cos \omega_0 t & 0 \leq t \leq t_0 \\ \left[-\frac{E_0}{\alpha_0} \int_0^{t-t_0} e^{-\alpha_0(t-t_0-\tau)} - e^{-\alpha_0(t-t_0+\tau)} \right] & t_0 \leq t \leq t_1 \end{cases} \quad (2-1a)$$

$$E(t) = \begin{cases} E_0 e^{-\alpha_0 t} \cos \omega_0 t & 0 \leq t \leq t_0 \\ \left[-\frac{E_0}{\alpha_0} \int_0^{t-t_0} e^{-\alpha_0(t-t_0-\tau)} - e^{-\alpha_0(t-t_0+\tau)} \right] & t_0 \leq t \leq t_1 \end{cases} \quad (2-1b)$$

where t_0 , t_1 , $t_1 - t_0$ are parameters related to the glottal flow peak, maximum closing rate and glottal closure, respectively. The parameter α_0 is used to control the steepness of source phase, and the parameter ω_0 , defined as $\omega_0 = 2\pi f_0$, determines the frequency of unvoiced. Parameters E_0 , α and ω are for computational use only. A typical LF model waveform is shown in Figure 2-3.

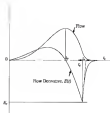


Figure 3-5 LF-model waveforms, dQ/ds for the differentiated ground flow

The first segment of the LF model characterizes the deflected glacial flow over the interval from the glacial opening to the maximum negative curvature of the waveform. The second segment represents a residual glacial flow that comes after the maximum negative-curvature. It can be shown from Eq. (2-1) that the spectrum of the first segment is dominated by the exponential component, e^{at} , of which the "negative bandwidth" equals $-a$. Likewise, the frequency response of the second segment can be approximated by a first order lag-pass filter with a cutoff frequency $F_2 = 1/(2\pi \tau_2)$ (Pien and Liu, 1963). As a result, the bandwidths of the first and second segments, B_1 and B_2 , are

$$B_1 = \frac{a}{2} \quad (2-2)$$

$$B_2 = \frac{1}{2\pi\tau_2} \quad (2-3)$$

It can be shown that the poles of the filter are either both real or complex-conjugate pair. The corner frequency ω_c and bandwidth B can be calculated from the zeros, z_i 's, by

$$\omega_c = \tan^{-1} \left[\frac{\operatorname{Im}(z_1)}{\operatorname{Re}(z_1)} \right] \quad (2-4)$$

$$B = -\ln |z_1| \quad (2-5)$$

We have found that the corner frequency ω_c of the poles and ω_0 of the LF model are nearly the same. Thus, we are only concerned with the change in bandwidth of the poles of the inverse filter. The bandwidth of source spectrum B_s is, in general, very close to B_1 , causing the waveform of the first segment to be undistorted after inverse filtering. However, B_2 is much higher than B . The second segment thereby retains its waveform after inverse filtering although the waveform phase may be different from the original. A typical example is given in Figure 3-4, which displays the spectra of the first and second segments of LF-model as well as the corresponding spectrum of the two-pole model. As a result, the waveform derived from the LF-model waveform has a flat spectrum envelope and exhibits a sharp peak at the conjunction between two segments, where the glacial closure occurs in the LF-model.

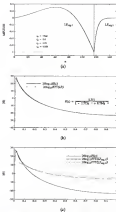


Figure 3-4: Effects of the inverse filter imposed on the differentiated glottal waveform: (a) LP model waveforms, (b) FFT spectra of LP model waveforms and two-pole model $N(1)$, (c) FFT spectra of individual segments of the LP model.

Knowing the relationship and transformation between the differentiated glottal and residue signals should enable us to extract one signal from the other. Although the residue signal does not appear to be highly informative, its integral, in contrast, tends to partially re-establish the shape of the first segment of the differentiated glottal flow. Thus, the analysis strategy based on the differentiated glottal flow can be transplanted to the integrated residue with little modification. To support our claim, we perform characteristic filtering based on a synthesis vessel, as shown in Figure 2-3, so that the similarities and contrasts between the LF-model waveform and the integrated residue can be noticed readily.

2.4/Choice of Model Type

In essence, building a speech model is equivalent to systematically combining the acoustic and perceptual attributes into a given construct. For speech synthesis, modeling is usually aimed at the parameterization of the voice source and the vocal tract.

In Chapter 1, we have shown that the speech production and propagation mechanisms can be described by a source-filter model (Fries, 1963), consisting of the glottal source, vocal tract transfer function, and lip radiation. This model is not only simple but effective in characterizing speech signals. The formant and LF synthesizers belong to this group and both synthesizers were widely used vocal synthesis (Kjarsgaard, 1973; Holman, 1973; Sambur et al., 1976; Aal and Ewert, 1978; Kawahara, 1984; Hornum et al., 1985; Holman, 1986; Kim, 1987; Miao et al., 1987; Childers et al., 1988; Childers and Wu, 1990; Kim and Kim, 1990; Childers and Lee, 1991; Lohman and Childers, 1991).

For the formant synthesizers, the properties of each component of the source-filter model are obtained individually. Usually, the lip radiation is approximated by a differentiator. The vocal tract transfer function is characterized by the formant anti-formants, which are implemented by resonance-resonance filters. The source model is designed to achieve the glottal volume-velocity waveform. Speech quality generated from

LP model



synthesized vessel \dot{V}



residue



integrated residue



Figure 2-3 Illustration of the similarity between the differential arterial flow and the integrated residue signal. Waveforms from top to bottom are: (1) LP model, (2) synthesized vessel \dot{V} produced by the LP model, (3) residue signal, and (4) integral of the residue signal.

such systematics was judged to be sufficiently high, provided that the global flow is appropriately modeled (Kahn, 1980; Hildner, 1983; Fain et al., 1988).

For the LP synthesizers, the composite frequency response of the global flow, vocal tract and lip radiators is modeled by a slowly time-varying filter (Audi and Hansen, 1971). The associated source excitation is primarily used to account for the production of the global pulse, in other words the pitch. The synthetic speech quality for many previous models is considered unsatisfactory due to an oversimplified excitation, failure to properly identify voicing, and poor spectral resolution (Ning, 1985; Kahn and Garst, 1985). However, with the use of sophisticated source excitations, the LP synthesizers can still achieve a very high quality.

On the whole, it appears that the perceptual quality of synthetic speech is improved by improving source excitation models for both LP and formant synthesizers. In the past, the lack of a physiological interpretation for the excitation was considered the primary obstacle against the use of LP techniques in generating a specific vocal quality, making the formant synthesizer a more popular tool. Nonetheless, this argument may no longer be valid once we are able to verify the relationship between the excitation and the global flow. Since the LP synthesizer has the advantages of (1) computational efficiency, and (2) ease of obtaining the excitation from speech, we decided to use the LP synthesizer model as the means to accomplish this research.

We start by integrating the acoustic stimulus into a comprehensive model. Our strategy in constructing a high-quality LP speech production model follows the analysis-by-synthesis idea. This speech production model is depicted in Figure 3-6, while Figure 3-3 presents the correlations we are going to examine between the acoustic parameters and model parameters. Before we work on the details, there is much groundwork to establish.

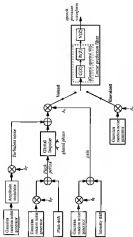


Figure 2-4. A, an individual from a population with a high frequency of the recessive allele.

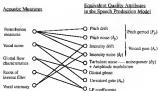


Figure 3.7: Correlations between acoustic attributes and model parameters.
 (The sentence "A → B" is read as "A is related to B".)

2.1 Data Collection and Methodological Considerations

In addition to the description of the experimental data base, this section provides background information about vocal quality including vocal fry and breathy voices. The acoustic properties quoted from previous research are listed for the purpose of comparing our analysis results. After explaining the measurement methods, we discuss the pre-processing schemes required to extract the acoustic features. These preparations constitute the foundation for the source extraction.

2.1.1 Experimental Data Base

The vowel *hi* was chosen as the experiment because it has been demonstrated for slow high speed laryngeal photography. Nine utterances for three different voice types, i.e., modal, vocal fry, and breathy, served as our data base. All three utterances were sampled by professional speech scientists. A description of the data base is shown in Table 2-1.

Table 2-1. Data base for speech analysis.

Subject	Sex	Voice type	# of pitch periods
S01	M	modal	303
S02	M	modal	303
S03	M	modal	264
S1	M	vocal fry	176
S2	M	vocal fry	209
S3	M	vocal fry	109
S4	M	breathy	280
S5	M	breathy	273
S6	M	breathy	454

During speech processing, the segments were performed over a steady state interval. As will be discussed later, the acoustic measures required a precise identification of the pitch period. An additional signal, the electroglottograph (EOG), was employed in this study to aid the speech processing. We sampled the speech and EOG signals at 10 KHz with 16-bit

precision. Each signal was digitized simultaneously using Digital Analog Converter DAC-048 preamplifier and ADC-100 digitizers. The microphone was an Electro-voice Model 88-12 held six inches from the lips. Before digitization, the signals were bandlimited to 5 kHz by anti-aliasing, passive, elliptic filters with a maximum stopband attenuation of ~ 35 dB and a passband ripple of ~ 0.7 dB. All data recordings were collected in an Industrial Acoustics Company (IAC) single-wall sound booth. To compensate for the microphone characteristics at low frequencies, the frequency response of the speech recordings was further corrected using a linear phase FIR filter.

2.1.2 Vowel Quality

The adequacy of an acoustic measure can be illustrated by an example for characterizing vowel quality. When assessing the acoustic measures, we certainly need to have a general concept of the vowel quality. As mentioned in the previous chapter, the vowel quality is referred to as the auditory response the listener experiences upon hearing the speech of another talker. Major types of vowel quality, according to Laver and Blount (1981), are modal, breathy, nasal lip, labiodental, hardness, and whisper. We excluded labiodental, hardness and whisper from this study because the other three vowel types were considered sufficiently representative of three modes of vocal fold vibratory patterns. The qualitative definitions (Lieberman and Blount, 1982; Lieberman et al., 1990) of the three vowel types are

Modal : Defined as a normal phonation. A modal phonation is characterized by a moderate frequency, wide laryngeal excursions, and complete closure of the glottis during about one third of the entire pitch period.

Breathy : Defined as audible escape of air through the glottis due to insufficient glottal closure. The degree of breathiness severity is inversely proportional to the length of the closed glottal phase.

Vowel fry is defined as a low-pitched, creaky kind of phonation. It also shows a great deal of irregularity from one pitch period to the next.

In this study, we are interested in ways the vowel characteristics change in the vibratory frequency of the laryngeal vibration. This is because the effect of the glottal vibration in terms of vocal registers is already reflected in the categorization of various voice types. Some acoustic features of glottal factors of various voice types are summarized in Table 2-2. These features will prove as relevant when we measure speech features using proposed acoustic measures.

Table 2-2. Summary of acoustic characteristics of glottal sources for three voice types.

		Modal	Vocal fry	Breathy
Fundamental frequency		medium	low	medium
Vibrational structure	Start	low	high	high
	Duration	low	low	high
Properties of glottal flow	Turbulent noise	medium	low	high
	Pulse width	medium	short	long
	Pulse duration	medium	high	low
	Asymmetry of closure	medium	fast	slow
Spectral tilt		medium	flat	steep
Vocal intensity		wide range	low	low

Source: Aho, 1994; Childers and Lee, 1991.

2.5.3. Analytical Logic

In most research the characteristics of glottal flow for one pitch period are determined using variables consisting of either the relative timing or the duration of special events such as the glottal opening and closure. Because the pitch period is usually a known value, it is

the waveforms, rather than the absolute timing, that has attracted the researchers' attention. This suggests a standardization procedure for those variables based on the underlying peak period. Properties of the standardized variables derived from a large population are assumed to represent general characteristics of the global system. Many problems and conclusions pertaining to macrodynamics are directly deduced based on the statistical results. Alternatively, such a statistical analysis can be performed by evaluating the averaged global pulse over a large number of sample periods. Such a logical analysis will facilitate the inquiry of some timing events in the global flow, the differentiated global flow and the integrated system.

1.5.1 Standardization of Peak Period

To perform the alternative statistical analysis suggested above, we resample every peak period at a variable rate so that every digitized waveform has the same length. In other words, the sampling rate for each individual peak period should be different in order to make the digitized waveforms symmetric. Difficulties associated with this problem lie in the identification and standardization of each peak period. A direct and exact solution, from a mathematical point of view, is the Sinc-resampled sampling rate conversion (Schulter and Roberts, 1971; Kruen and Arai, 1990; Schreiner and Chirp, 1990). The Sinc interpolation is given by

$$\hat{x}(xT) = \sum_{n=-\infty}^{\infty} x(n) \frac{\sin\left\{x_0(xT - \frac{1}{2} + n)\right\}}{x_0(xT - \frac{1}{2} + n)} \quad (1-4)$$

where x_0 is the sampling frequency of the original sequence $x(n)$, T is a new sampling interval of $\hat{x}(x)$, $\sin(\theta)$ is a phase offset. The T and θ for each individual period are determined so as to yield a maximum similarity across the resampled periods.

The promising task required by Sinc interpolation is computationally expensive. A simpler approach is presented below to facilitate the computations. First of all, we interpolate

the analyzed signal $x(n)$ by a factor of five times by using a lowpass filter:

$$\hat{x}(n) = x(n) * \hat{h}(n) \quad (2-7)$$

where $*$ denotes the convolution, $\hat{x}(n)$ is the linearly interpolated data sequence, and $\hat{h}(n)$ is the impulse response of a lowpass filter with the cut off frequency at $\pi/5$. In our case, a 511-order FIR filter designed by using the window method is employed to avoid phase distortion. The impulse response of this FIR filter is

$$\hat{h}(n) = \frac{1}{L} h(n/L) \quad (2-8)$$

$$\text{where: } h(n) = \frac{\sin(\pi n \frac{1}{5})}{\pi n \frac{1}{5}} \quad \text{for } n = 0, \pm 1, \pm 2, \dots \quad (2-9)$$

$$w(n) = \begin{cases} 34 - 46 \cos \frac{2\pi n}{103} - 1, & \text{if } n \leq 103 \\ 0, & \text{otherwise,} \end{cases} \quad (2-10)$$

$$L = \sum_{n=-103}^{103} h(n)w(n) \quad (2-11)$$

The next step is to separate each individual pitch period along the signal. We use the two-channel approach (Kochenevskiy and Chelstov, 1989) to obtain this processing automatically. The glottal closure instant, which is signaled by a rapid decrease in the DCO, has been found to coincide well the maximum in the differentiated DCO (DDCO) for that period. Thus, we can locate the instant of glottal closure by picking the negative peaks of the DDCO signal, as illustrated in Figure 2-6. A pitch period is then defined as the interval between two consecutive glottal closure instants.

Due to the propagation delay of the sound wave from the glottis to the microphone, we apply a time lag of 0.5 msec to the DCO signal to achieve synchronization with the speech signal. Also, in order to improve the accuracy of the locations of the peaks, we employ a quadratic interpolation method (Markel and Gray, 1976; Tsao et al., 1987). Let $P = (1, P(0))$,

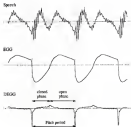


Figure 3-8: Synchronized speeds, RGO, DEGO signals.

and $f(1)$ define three points centered at $g(k)$, where $f(0)$ corresponds to a discrete minimum value, and $f(-1)$ and $f(1)$ are points to the left and right of $f(0)$. The position of interpolated maximum, k , is then resolved from a second-order approximation using three points by

$$k = -\frac{f'(0)}{f''(0)} = -\frac{f(1) - f(-1)}{2f(1) - 4f(0) + f(-1)} \quad (2-13)$$

Values obtained through this process are rounded to the nearest sample of the re-sampled signal. Thus, the resulting resolution of each pitch period, estimated from the re-sampled response, increases by approximately three times.

Finally the length of the pitch period is adjusted to 512 samples by using the FFT method. That is, depending on the number of the samples in the interpolated period, we append zeros or remove the high frequency range of the FFT sequence to achieve the intended length (512 samples in this case). The fixed-length signal is then obtained by taking the IFFT of the modified FFT sequence. Notice that the discontinuity (linear trend) between two boundaries of the underlying signal must be removed before applying the FFT method since the signal has to be circularly periodic.

2.3 Feature Extraction

Under the guidance of the proposed speech production model (Figure 2-6), we explore the relations between the model parameters and some relating acoustical measures. A pitch-synchronous covariance LP analysis is adopted to measure the spectral properties of the speech signal. The LP order is chosen to be 14 to account for the spectral rich of the glottal flow and the number of formants within 5 kHz bandwidth. Following the approaches presented in the survey of acoustic measures, we illustrate how to extract model parameters that correspond to the perturbation measures, spectral info, phase characteristics, and vocal source sequentially. The relationships regarding the extraction of acoustic features are considered in this study.

3.6.1 Perturbation Measures

As mentioned previously, the perturbation measures we used to characterize the vibratory patterns of the vocal folds, which include variations in the pitch period and waveform amplitude. The perturbation measures can be further divided into two types, namely, deterministic and random noise as demonstrated in Figure 3-9. The subharmonics result from a repetitive vibratory pattern extending more than one-pitch period, while the random noise represents the unpredictable characteristic of the vocal fold vibration. To avoid further complicating the problem, we confined our research to random noise while discussing the proposed perturbation measures.

Because each subject involved in this experiment was instructed to utter a steady vowel with a comfortable intensity, the pitch and intensity contours of recorded speech were considered to be fairly stable. Typical pitch and intensity contours of a model voice are shown in Figure 3-10(a) and Figure 3-11(a). As we say mentioned in the perturbation associated with the measured signal, a proper initial step is to obtain the corresponding deviation by removing the average value. By inspecting the spectral properties of the deviations of pitch and intensity signals (Figures 3-10(b) and 3-11(b)), we find that both signals are relatively flat except in the region of low frequencies. This finding leads us to conjecture that the deviation signal can be modeled as a slow fluctuating component accompanied with a white noise source. The low-frequency component in the deviation signal, termed "drift" as many studies, is known as the coherent nature of human speech. Though the dynamic patterns of drift determine the start of speech, they are unlikely to provide much information about vocal quality. Thus, we set out to characterizing this effect while studying vocal quality. In fact, as pictured from the point of view of the Elzinga process, most perturbation measures were introduced to characterize the drift or to compensate for white noise source. Examples for some perturbation measures and their mathematical relationships can be seen in Pardo and Thue (1992).

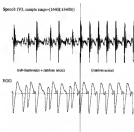


Figure 2-6: Demonstration of two types of perturbations: with-noise and random noise.

The use of a high-pass filter will not properly separate the noise source from the drift because it removes the low-frequency portion of the noise also. Thus, we performed the separation in the frequency domain using a DFT method with the following steps:

1. We remove the linear trend of the test and boundaries of the sampling signal to avoid large discontinuities occurring at boundaries due to the DFT method.
2. Before computing the DFT sequence of the residual deviation signal, we further remove the d.c. component introduced by the first step.
3. Except the d.c. component, the magnitude of the DFT sequence below the 4th sampling rate (400) is set to the average magnitude of the rest DFT sequences (see Figures 2-10(a) and 2-11(b)).
4. Over the new DFT sequence with frequency unchanged, we then take the inverse DFT of the new sequence to yield the same signals (see Figures 2-10(c) and 2-11(c)).

Examples of the histograms of the noise signals are shown in Figure 2-12. It appears that the zero-mean Gaussian distribution provides a good fit for the underlying frequency distribution. We therefore assumed that the noise component satisfies a Gaussian distribution, in which the standard deviation is sufficient to characterize the statistical property. This hypothesis can be informally validated by comparing the cumulative probability density functions of the noise components and by comparing them with the corresponding Gaussian distribution function with the same mean and variance (Figure 2-13).

Following the nomenclatures defined by Paoletti and Titter (1983), we use δ_p and δ_s to denote the standard deviations of the pitch and intensity noise components respectively. The following discussion characterizes how the δ_p and δ_s are related to the jitter and shimmer. We define the normalized jitter (in percent) as

$$\% \text{ jitter} = \frac{1}{L} \sum_{i=1}^L \frac{|P_i - P_{i-1}|}{P_i} \times 100\% \quad (2-17)$$

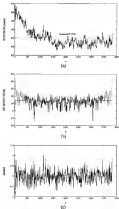


Figure 3-18. Normalized of peak noise: (a) original peak noise; (b) magnitude DFT (dashed line) of the deviation signal and the one after adjustment (solid line); (c) peak noise after taking the inverse DFT of the adjusted DFT response.

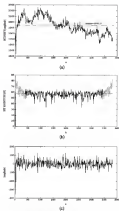


Figure 2-11: Extraction of monthly noise: (a) intensity contour; (b) and (c) are the same as in Figure 2-10.

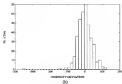
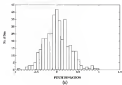
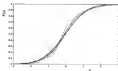
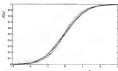


Figure 3-12. Histograms of (a) path noise, and (b) intensity noise.



(a)



(b)

Figure 3-43 Normalized cumulative probability distribution functions $[F(x)]$ of perturbation noise for noise structure: (a) peak noise, (b) narrow noise (shown as dotted lines). The corresponding Gaussian distribution with the same variance is shown by the solid line.

where P_0 denotes the absolute value, P_1 is the n th pitch period in a segment of n pitch periods and P_0 is the mean value of P_1 's. If we define $P_d^1 = P^1 - P_0$, then Eq. (2-13) can be rewritten as

$$\begin{aligned}\% \text{ jitter} &= \frac{1}{P_0} \sum_{i=1}^n \frac{P_1 + P_d^1 - (P_0 + P_d^{1-1})}{P_1} \times 100\% \\ &= \frac{1}{P_0} \sum_{i=1}^n \frac{P_d^1 - P_d^{1-1}}{P_1} \times 100\% \quad (2-14)\end{aligned}$$

If we assume P_d to be a random process with a zero-mean Gaussian distribution and a var Σ , the equation can be approximated by

$$\begin{aligned}\% \text{ jitter} &= \frac{\sqrt{\Sigma}}{P_0} \sqrt{\frac{1}{n}} \times 100\% \\ &= \frac{\sqrt{\Sigma}}{P_0} \sqrt{\frac{1}{n}} A_p \times 100\% = \frac{A_p}{\sqrt{n}} A_p \frac{100\%}{P_0} \quad (2-15)\end{aligned}$$

where the number denotes the statistical mean. In a similar manner, the percent standard can be derived as

$$\begin{aligned}\% \text{ standard} &= \frac{1}{P_0} \sum_{i=1}^n \frac{(P_1 - A_p)_{-1}^2}{A_p} \times 100\% \\ &= \frac{A_p}{\sqrt{n}} A_p \frac{100\%}{A_p} \quad (2-16)\end{aligned}$$

where A_p denotes the square root of the average time (power) of the n th glottal period, and A_p is the average of A_p 's. Compared to the definition given by other researchers, where A_p is the peak magnitude, the adopted form is more likely to correspond with the perceptual characteristics of the form as laboratory experiments that involve short-term spectra of acoustic signals. More important, it is the power density rather than the peak magnitude used for the speech analysis and synthesis in the proposed speech production model.

From Eqs (2-12) and (2-14) we know that the \hat{W}_{ij} and \hat{V}_{ik} denote two particular types of mean forecasts for the pack and inventory series. To verify the forecasting effectiveness, we list the \hat{A}_j 's and \hat{B}_j 's measured from the pack and inventory series signals as well as those derived from \hat{W}_{ij} and \hat{V}_{ik} in Table 2-3. The obtained values are very close to the values computed from two different approaches, each under the very substantial assumptions with regard to the perturbation value and its relation to other perturbation measures.

Table 2-3. Mean values and standard deviations (in italics) of the accurate measures for each subject.

	True period (ms)	\hat{W}_{ij}	\hat{A}_j (in ms)	\hat{B}_j	\hat{A}_j (in ms)	\hat{V}_{ik} (ms)	\hat{A}_j (in ms)	\hat{B}_j
M1	2703 <i>±0.103</i>	0.381	0.294	0.299	1.000 <i>±0.072</i>	2.352	43.473	43.822
M2	7731 <i>±0.049</i>	0.377	0.237	0.232	1.000 <i>±0.118</i>	3.787	27.341	24.663
M3	7709 <i>±0.062</i>	0.344	0.235	0.249	1.000 <i>±0.053</i>	2.763	89.334	72.297
V1	11183 <i>±0.089</i>	0.366	0.337	0.315	1.000 <i>±0.059</i>	3.408	44.734	43.703
V2	7590 <i>±0.126</i>	0.434	0.337	0.297	1.000 <i>±0.233</i>	2.754	44.387	41.408
V3	38763 <i>±0.723</i>	0.383	0.263	0.243	1.000 <i>±0.348</i>	4.043	343.521	60.608
B1	9794 <i>±0.193</i>	0.693	0.733	0.746	1.000 <i>±0.053</i>	2.755	230.283	233.234
B2	9768 <i>±0.223</i>	0.603	0.634	0.769	1.000 <i>±0.104</i>	2.143	860.295	1187.19
B3	2487 <i>±0.043</i>	0.473	0.383	0.377	1.000 <i>±0.084</i>	0.847	91.799	84.949

2.4.2 Spectral tilt

Theoretically, the transfer function of the vocal tract is characterized by a set of formants, which are distributed along the frequency axis. The spectral tilt of speech is mainly dominated by the lip radiation and glottal shaping filter (Rabiner and Schafer, 1978). If the lip radiation is modeled as a differentiator, then the differentiated glottal flow becomes the most pertinent component to determine the spectral tilt of a speech signal. This implies that the spectral tilt of the differentiated glottal pulse can be estimated from the speech signal.

As discussed in Section 2.1, we model a two-pole filter to approximate the spectral tilt of differentiated glottal flow. The filter coefficients are then estimated using LP analysis based on the speech signal. Table 2-5 lists the estimated LP coefficients for the three different voice types.

2.4.3 Glottal phase characteristics

In addition to a general comparison of the glottal phase characteristics for the three utterances, we present a novel measure called “integrated residue” to depict the entire phase of the glottal flow.

2.4.3.1 General properties

Because the magnitude spectrum of the residue signal is flat (not in a strict sense if we consider the spectral harmonics and modeling errors) due to the inverse filtering, phase characteristics are the only information left in the residue that is related to the glottal source (Wang and Maekel, 1978; Houton, 1983). As mentioned earlier, it is the integrated residue which resembles the differentiated glottal flow, reflecting pertinent physiological features of the vocal folds. Presumably, we can explore the phase properties of the glottal source by examining the integrated residue. Unfortunately, the timing of unvoiced glottal events are distorted by the inverse filtering and integration. Those features extracted from the integrated residue are not as useful as those for the glottal flow and, therefore, will be not

unweighted. Instead, the comparison across the integrated modulus of the main vibrations is performed using correlation coefficients. The results are given in Table 3-4. These results, however, do not adequately characterize the correlation between the glottal phase and vocal quality. Consequently, we introduce another measure called the “abruptness index” to measure the information about the return phase of the glottal flow.

Table 3-4. Correlation coefficients across all vibrations.

	B0	B2	B1	V3	V2	V1	M3	M2	M1
B0	1.0000	0.2330	0.7437	-0.5287	0.4437	0.3804	0.8607	0.9437	1.0000
B2	0.2330	1.0000	0.8006	-0.4217	0.9761	0.3608	0.9086	1.0000	
B1	0.7437	0.8006	1.0000	-0.1909	0.3381	0.2379	1.0000		
V3	-0.5287	-0.4217	-0.1909	1.0000					
V2	0.4437	0.9761	0.3381	-0.4217	1.0000				
V1	0.3804	0.3608	0.2379			1.0000			
M3	0.8607	0.9086	1.0000				1.0000		
M2	0.9437	1.0000						1.0000	
M1	1.0000								1.0000

2.6.3.2. Abruptness index

The idea for this measure comes from the LF-model. In a study of the acoustic variability of the glottal source factors, Adas [1991] concluded that the q_1 , q_2 of the LF-model were the two most significant parameters contributing to vocal quality. Between q_1 and q_2 are the parameters controlling the opening of the return phase in the LF-model, they are accepted as an indicator of vocal abruptness. In Adas's study, q_1 was defined as the instant at which the amplitude of the modulated differentiated glottal flow dropped to 1% of its peak value.

Accordingly, \hat{z}_1 can be derived by solving the following equation

$$\frac{e^{-j(\hat{\omega}_0 - \Omega)} - e^{-j(\hat{\omega}_0 - \Omega_0)}}{1 - e^{-j(\hat{\omega}_0 - \Omega)}} = \hat{z}_1 \hat{z}_2 \quad (2-17)$$

where \hat{z}_2 is the pitch period and \hat{z}_1 can be obtained a priori by solving

$$\hat{z}_{10} = 1 - e^{-j(\hat{\omega}_0 - \Omega_0)} \quad (2-18)$$

Since \hat{z}_1 and \hat{z}_2 form a mathematical mapping, we may say that the statistical significance of these two parameters stand on the same footing. Based on this understanding, we need to focus on only one parameter. It can be shown from Eq. (2-1) that

$$\frac{dR(\hat{z})}{d\hat{z}} = \frac{E}{\hat{z}} \quad \text{or} \quad \hat{z}_0 = E_{\max} \left(\frac{dR(\hat{z})}{d\hat{z}} \right) \quad (2-19)$$

The equations above explicitly tell us that \hat{z}_0 can be readily obtained if we know the derivatives of $R(\hat{z})$ at \hat{z}_0 . Since the constant E_{\max} usually coincides with the largest value of $dR(\hat{z})$ and E_{\min} is usually the minimum of $R(\hat{z})$, it will be convenient for us to calculate \hat{z}_0 using the following equation

$$\hat{z}_0 = \frac{-\max(dR(\hat{z}))}{\max(dR(\hat{z}))} \Delta t \quad (2-20)$$

where Δt denotes an infinitesimal time interval. For a discrete signal, this value, Δt , can be substituted by the sampling interval, ΔT , provided that the interval is sufficiently small. If we define the vocal sharpness index, \hat{z}_0 , as the normalized \hat{z}_0 in percentage, then \hat{z}_0 becomes

$$\hat{z}_0 = \frac{-\max(dR(\hat{z}\Delta T))}{\max(dR(\hat{z}\Delta T))} \frac{\Delta T}{P_0} \times 100\% \quad (2-21)$$

where ΔT is the sampling rate, P_0 is the pitch period, and $dR(\hat{z})$ stands for the difference function. It is obvious that \hat{z}_0 is readily obtained once $R(\hat{z}\Delta T)$ is available.

The importance of the $\hat{U}_2(\omega; T)$ calls for an employment of the phasor structure filtering technique, which is not always feasible. Thus, we demonstrate here to convert the modulator signal into the differentiated phasor flow. As discussed in Section 2.3, a two-pole filter was employed to model the spectrum slope of the differentiated phasor flow $\hat{U}_2'(\omega; T)$. The transfer function of $\hat{U}_2'(\omega; T)$ is given as

$$\hat{U}_2'(\omega) = \frac{\omega_0^2}{1 - \omega_0^2 \omega^{-2} - \omega_0^2 \omega^2} \quad (2-22)$$

where ω_0 denotes the 2-resonance of the modulator. The phasor differentiated flow, $\hat{U}_2'(\omega; T)$, can be approximated by fixing the modulator signal into the two-pole filter, of which the coefficients is derived by LTP analysis of the speech signal. Substituting $\hat{U}_2'(\omega; T)$ for $\hat{U}_2(\omega; T)$ in Eq. (2-21), we can easily obtain the sharpness index (listed in Table 2-5 on page 63).

2.4.4 Vocal Noise

Much research has been devoted to estimating vocal noise performance in a variety of situations or training speech. However, due to the limited capability of existing technology, the noise could only be presented in a form of signal-to-noise or harmonic-to-noise ratios that does not offer enough details of vocal noise. To gain a better understanding of vocal noise, we plan to examine the properties of noise from three aspects, namely, the signal-to-noise ratio (SNR), amplitude modulation, and noise spectra.

2.4.4.1 Noise extraction

In order to acquire the vocal noise techniques for identifying and separating the prototype patterns (i.e., the standardized signal of one pitch period along an utterance) we required. The extraction of such periods was accomplished by peak picking the OGG signal. The separation of noise from the prototype can be achieved either in the frequency domain or in the time domain. There were two approaches that influenced us most in carrying

to accomplish the noise extraction. One of the approaches was proposed by Yumoto et al. (1982), who considered the noise as the deviation of a quasi-periodic speech signal. They first identified a prototype period of phonation by averaging the waveforms of every period in a steady utterance. This prototype was then subtracted from the speech signal for each pitch period to yield the noise. The use of such an approach, however, requires that the subject's utterances have to be strictly steady for a number of periods. This may not be feasible in some cases. Thus, Kameya et al. (1984a,b) proposed another measure by estimating the harmonic-to-noise ratio. The noise signal was isolated from periodic components either by using a comb filter or by reflecting the anti-harmonics in the spectrum of the analyzed signal. In such a method, any component not harmonically related to the fundamental frequency was classified as noise.

Although Kameya's method was robust and efficient in computation, it was inappropriate from the practical point of view since it took account of the jitter and shimmer. In our proposed speech production model the vocal noise is considered to be an independent module. This design concept requires that the uncontrolled factors, such as jitter and shimmer, has to be segregated from the real noise source. Thus, we adopt Yumoto's idea in a modified form. First, we standardize the power and length of every pitch period of the synthesized vowel using the method discussed in Section 2.5.4 before extracting the prototype. Then we calculate the vocal noise by minimizing the least square difference, $E(\hat{X}_i, \hat{X}_p)$, between the prototype signal \hat{X}_p and the analyzing signal \hat{X}_i :

$$E(\hat{X}_i, \hat{X}_p) = \sum_{n=0}^{N-1} |K^C(\hat{X}_i(n)) - \gamma \hat{X}_p(n)|^2 \quad (3-11)$$

where C^2 denotes the circular shift with an l lag, and N is the length of the standardized pitch period. The purpose of using C^2 is to rectify the phase discrepancy between \hat{X}_i and \hat{X}_p . The scale factor γ is then determined by using:

$$\partial E S_{\gamma}^2 S_{\gamma\gamma} / \partial \gamma = 0, \quad (2-14)$$

which leads to

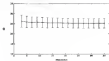
$$\gamma = \frac{\sum_{k=0}^{N-1} C^*(S_k) S_k(t_0)}{\left(\sum_{k=0}^{N-1} S_k^2(t_0) \sum_{k=0}^{N-1} S_k^2(t_0) \right)^{1/2}}. \quad (2-15)$$

The lag γ that yields the maximum γ is chosen as the cross-correlation offset. Thus, the noise signal becomes

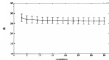
$$w(t) = C^*(S)(t) - \gamma S_{\gamma}(t). \quad (2-16)$$

A problem of this approach above is related to decreasing an appropriate number of periods to form an analysis window. The prototype derived from a short window is statistically unreliable. On the other hand, it is unlikely that a steady phenomenon is maintained throughout the stimulus. Therefore, we have to examine the influence of the number of periods with regard to the noise measure. An empirical but sensible criterion for selecting the analysis window is to search for the minimum period that gives a small standard deviation. Here we use three different stimulus as signal experiments. As shown in Figure 2-14, the standard deviations of these samples are relatively large when the analysis window is small. When the analysis window is increased to more than 13 periods, both the standard deviations and mean values become stable. Thus, the prototype is calculated using a window of 13 consecutive periods, in which the current period is located at the center of the window. The selected number is somewhat smaller than that suggested by other researchers (Ninio et al., 1982; Hsu et al., 1987; Eklund et al., 1990). We reason that this result is due to the nonlinear enhancement and peak randomization.

(a) model (B1)



(b) model (B2)



(c) model (B3)

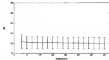


Figure 2-34: Variance of B1R versus number of analyzed periods
(range of each series test = [-0.01, and 0]).

Another problem arising concerns the detection of low frequency power. Because the air sampling from the lungs is not controlled, a frequent change of low frequency components is anticipated. Fortunately, owing to the fact that the noise in the low frequency region is perceptually masked by the harmonics of the fundamental frequency (Cholten and Lee, 1994), we can apply a notch filter to eliminate the low frequency components without disturbing the perceived quality. The cut-off frequency of the highpass zero-phase filter can be designed to adapt to the current pitch period P_i such that low-frequency components below 300 Hz are sufficiently suppressed and high-frequency components are not affected. The frequency response of the highpass filter is given by

$$H(z) = \frac{2 - \alpha}{2} \frac{1 - z^{-1}}{1 - \alpha z^{N-1}} \quad (3-27)$$

where $\alpha = \frac{(0.2P)}{322}$ provides the adaptation for the i th pitch period of length P_i . The number 312, is the length of the pitch period after oversampling. The scale factor, $(2-\alpha)/2$, in Eq. (3-27) is to make the magnitude unity at the one-half the sampling frequency. Eventually, the desired noise signal is the result after passing $x(i)$ through the notch filter.

3.6.4.2 Properties of vocal noise

Once we get the desired noise, the properties can be evaluated as:

- **Signal-to-Noise ratio**

The Signal-to-Noise Ratio (SNR) for the i th pitch period is calculated as

$$SNR_i = 10 \log_{10} \left[\frac{P^2 \sum_{n=1}^N x^2(n)}{\sum_{n=1}^N s_n^2(n)} \right] \quad (3-28)$$

and the SNR's for the noise situations are listed in Table 3-3

Table 1: 3. Mean values and standard deviations (\pm std) of the selected measures for each subject.

	α_1	α_2	β_0	$3\sigma\theta$ [deg]
M1	-0.581 ± 0.043	-0.179 ± 0.018	0.893 ± 0.165	25.466 ± 2.499
M2	-0.718 ± 0.037	-0.006 ± 0.063	1.390 ± 0.189	20.886 ± 2.138
M3	-0.757 ± 0.054	-0.159 ± 0.028	1.448 ± 0.155	25.329 ± 1.795
V1	-0.646 ± 0.057	-0.079 ± 0.051	0.681 ± 0.137	36.468 ± 1.538
V2	-0.584 ± 0.089	-0.092 ± 0.067	0.858 ± 0.291	24.822 ± 2.945
V3	-0.210 ± 0.157	0.177 ± 0.090	0.583 ± 0.090	18.768 ± 3.404
B1	-0.973 ± 0.080	-0.058 ± 0.067	1.371 ± 0.284	19.382 ± 2.945
B2	-1.361 ± 0.150	0.272 ± 0.148	2.645 ± 1.167	6.206 ± 2.336
B3	-1.657 ± 0.093	0.602 ± 0.091	5.135 ± 1.229	14.864 ± 2.705

• Amplitude modulation

As shown in Figure 2–15, the amplitude modulation is obtained by averaging the magnitude of the noise signal over all periods.

• Noise spectrum

The spectrum of the noise signal is computed using the FFT. Though the length of every pitch period has been expanded to 312 samples, the frequency resolution of the FFT sequence still depends on the actual fundamental frequency. Due to the fact that the fundamental frequency may change from period to period, the resolution of each FFT sequence is therefore different from each other. Thus, we apply the bilinear interpolation on the FFT sequence to achieve a unique frequency resolution. The individual FFT spectra are then averaged to yield an estimation of the noise spectrum. Figure 2–16 shows the period spectra for different voice types.

2.5.4.3 Final summary

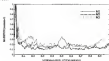
To summarize, we present the noise reduction algorithm as a flowchart in Figure 2–17. We recall that the noise is estimated from the residual signal, which is obtained using techniques addressed in previous sections. The overall procedure is tedious, but is straightforward and easy to implement.

2.2 Discussion

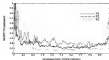
The results we have gained so far can be summarized as four aspects:

COF was found that the perturbation noise can be modeled by a zero-mean Gaussian process with a low-frequency drift. Measures that sufficiently downplayed the drift could be used to characterize the noise perturbations. In particular, we have used the 50-year and 70-decades to indicate the standard deviation of the noise sources. The results of measured perturbations with respect to three voice types, in general, were consistent with other researchers' reports, i.e., vocal fry and breathy voices exhibit higher perturbations. We also

(a) model



(b) vocal fry



(c) breathy

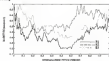
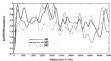
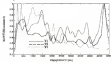


Figure 2-18. Amplitude modulation of words used for different vocal types.

(a) modal



(b) vocal fry



(c) breathy

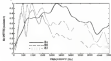


Figure 3-18. Spectra of vocal tones for different voice types.

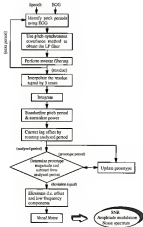


Figure 3-17. Schematic flowchart for noise correction.

found that the smaller perturbations in vocal fry and breathy voices corresponded to low pitch subjects. Interpreted from a psychoacoustic perspective, such values would have different impacts to the perception of vocal quality (Blissfield, 1937). Furthermore, loudness, a perceptual descriptor of the intensity, was reported to be a nonlinear function for various frequencies (Robinson and Chaffin, 1976). These factors confused the study of vocal quality mostly on the basis of quantitative measures. To achieve a thorough understanding of vocal quality, the research scope should cover speech perception as well (Plassaun, 1971a; Blissfield and Lanchion, 1944; Horrmansky et al., 1985; Wang et al., 1981).

(I) A comparison of the spectral tilt of the source can be performed by visually inspecting the frequency responses of the two-pole filter model. As shown in Figure 3-18, the spectral tilt is moderate, relative flat, and steep for vocal fry, modal, and breathy voices, respectively. A simple quantitative measure can be achieved by comparing the first coefficient a_1 , since the coefficients a_1 and the poles $z_{1,2}$ of the modeled filter have the following relation:

$$\begin{aligned} a_1 &= -2\Re\{z_{1,2}\} = -2\Re\{z\} & \text{if } \theta \neq 0, \pi \\ &= -|z_1| + |z_2| & \text{if } \theta = 0, \pi \end{aligned} \quad (3-29)$$

where $\theta = \tan^{-1}\left\{\frac{\Re\{z_1\}z_2^*}{\Re\{z_1\}^2}\right\}$

Herein, the value of a_1 can be used as indicator here since the poles are to the unit circle. A larger a_1 corresponds to a flatter spectral tilt and broader bandwidth. According to the data in Table 3-3, the values of a_1 for different voice types satisfied the following inequality:

$$|a_1(\text{vocal fry})| > |a_1(\text{modal})| > |a_1(\text{breathy})|$$

This result is consistent with the previous observation in Figure 3-18 and with the conclusion shown in Table 3-2.

(2) As the relation between the nasals signal and glottal flow was revealed, the glottal phase characteristics could be traced back from the reaction signal. We have studied

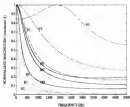


Figure 3-18. Frequency response of two-pole filter for nine subjects

the similarities and differences among the nine integrated vowel signals by examining their correlation coefficients. It was found that no similarity exists among the model vowels. This result suggests that the glottal phase characteristics cannot be described with a general pattern. On the contrary, the shapeless index showed an advantage for characterizing vowel types. The values of the shapeless indices for these vowel types are presented in descending order as modal, fry, modal, and breathy vowels. Such a measure may enable researchers to classify vowel types with considerable convenience. It is worth noting that the meaning of this measure can be interpreted from various aspects. In the time domain, it is related to the temporal transition from the maximum glottal closure instant to the glottal opening. In the frequency domain, it indicates the spectral slope of the glottal source. From the point of view of the neural signal, it corresponds to the peak factor of the main excitation pulse.

(4) The measured SHR^2 's for various vowel types supported the earlier finding of other researchers that breathy vowels, in general, were accompanied by the largest vocal noise. Such a noise level was effective enough to undermine its role in source modeling. An intriguing observation with the breathy vowels is that the standard deviations of corresponding SHR^2 's are as small as that of modal vowels. This result suggests that a steady noise source would be appropriate to model the vocal noise for modal and breathy vowels.

As displayed in Figures 2-14, the noise spectra for different phonemes were fairly flat, suggesting that the noise for the integrated vowel is white. However, for the purpose of speech synthesis, the noise source has to be pre-emphasized by a Highpass Filter before applying to a LP synthesizer.

The amplitude modulation of the vocal noise generally mimics the magnitude of the integrated vowel (Figure 2-15). However, the high amplitude modulation near the glottal closure may also be related to the phase misalignment. Notice that there are two types of modulations presented in the modal signal: one is the speech variation caused by the vocal fold closure, and the other is the variation due to the airflow turbulency from the lungs (Kang and Ewert, 1983). The integration with respect to the random makes the phase

alignment procedure in favour of the airflow, thus increasing the degree of mismatching for the speech variation. The assumption above can be further verified by the subsequent derivation. Suppose there is a phase offset, θ , between two signals, $x(t)$ and $x_1(t) = \theta(t)x(t - \tau)$, then the difference $e(t)$ is

$$\begin{aligned} e(t) &= x(t) - x_1(t) \\ &= x(t) - \theta(t)x(t - \tau) \\ &= 2x(t)\sin^2\frac{\theta(t)}{2} \\ &= x(t)\frac{\theta(t)}{2}, \quad \text{for } \theta(t) \ll \frac{\pi}{2} \end{aligned} \quad (2-38)$$

Clearly from Eq. (2-38), the error $e(t)$ is proportional to the signal, resulting in the consistency between the amplitude modulation of the reconstructed and the magnitude of the original signal.

So far it is undetermined whether the amplitude modulation is an artifact of the analysis method or is a primitive feature of the glottal source. We will reveal this issue in Chapter 4. But one thing is certain here, that is, the quality of synthetic speech is affected by the modulation of vocal tract.

2.6 Conclusion

In this chapter we have explored the acoustic features within the source-filter theory. The properties of the glottal source were primarily extracted from the integrated residual signal, which was obtained by making use of the pitch-synchronous LP analysis with the aid of the DDDG signal. We demonstrated the analysis methods using sustained vowels, /a/, of three voice types, i.e., modal, vocal fry, and breathy voices. The roles of many existing acoustic measures were carefully investigated. Although more extensive investigations are needed in order to establish statistical significance of model parameters, the results of our study provided a basic understanding of source variations as well as their manifestations in the acoustic measures. More important, the capabilities of extracting the

glottal source properties using LP analysis were substantiated. The competence of LP method in speech analysis suggests that a high quality LPC synthesizer is achievable. Of course, this is under the assumption that the properties of the glottal source are faithfully preserved during the analysis and synthesis. To achieve such a requirement, two important features that are usually ignored in many LP synthesizers, i.e. the vocal tract and the glottal phase characteristics, have to be incorporated into the source model.

CHAPTER 3 SOURCE MODELING

Since it was introduced in the late 1960s, the linear predictive coding (LPC) technique has been extensively used in speech processing and coding (Rabinowitz and Schuster, 1979). Speech systems are considered under class of LPC coders as a slowly time-varying all-pole filter to model the composite spectral characteristics of the glottal flow, vocal tract and lip radiation. The excitation for this all-pole filter is a sequence of signal with quasi-periodic phases for voiced speech and random phases for unvoiced speech. In this study, we apply a sixth order polynomial model to determine the phase characteristics of voiced source excitation. Source filters derived by this model are further compressed through a vector quantization technique. A 33-entry glottal codebook is derived by quantizing the voiced samples stored by 33 codewords. On the other hand, a 256-entry stochastic codebook is generated for unvoiced speech systems. However, unlike the glottal codebook, codewords in the stochastic codebook are simply taken from a Gaussian noise source.

3.1 Review of Previous Research

Over the years, various types of excitation have been proposed to drive the synthesis filter to produce speech. In the conventional pitch-excited LPC vocoder (Noll and Huxman, 1971), the excitation signal is either an impulse train for voiced speech or a random noise for unvoiced speech. The quality of synthesized speech in some applications is judged as unnatural due to incorrect timing elements, poor spectral resolution and oversimplified radiation functions (Wong, 1980; Rabin and Gersu, 1983).

The use of either synthesized-excitation functions such as the Multi-Pulse (MP), Code-Excited (CE) or their variants (Arai and Ikeda, 1982; Schroeder and Arai, 1985; Wang and Arai, 1989; Koss and Burnswell, 1990) can result in high-quality synthetic speech if the synthetic excitation is described sufficiently well by adequate number of codewords or pulses. Coders using this type of excitation go beyond spectral analysis and pitch estimation. Parameters representable by predictive filters can be recovered by formulating the excitation signal. That is, the excitation signal is formed by searching for the best candidate in a given set of candidate sequences by minimizing the spectrally weighted difference between the original and the synthesized speech signals.

In fact, the ideal excitation for LP synthesizers is the residue signal obtained by inverse filtering of the original speech signal. Attempts have been made to encode and transmit the residue signal in many coding systems (Lin and Mapp, 1975; Choudhury and Wang, 1979). But little research effort has been directed to extracting the features of the residue signal. In 1978 Wang and Mapp's synthesized prototype excitation pulse by inverse filtering the deconvoluted glottal flow of the vowel /a/. Although this excitation pulse was intentionally designed to reduce the features of synthesized speech, both quality and naturalness, as expected, were improved due to the preserved glottal characteristics. However, the excitation pulse presented in their experiment has certain drawbacks. First, it is feasible only when the fundamental frequency is below 100 Hz. Second, a single prototype excitation pulse is not likely to suit all the situations since glottal features for various speakers and phonemes can vary considerably.

The importance of glottal characteristics for speech synthesis was also demonstrated by Ferguson and Flanagan (1985). Finding the similarity between the residue and the second derivative of the glottal pulse, they incorporated the glottal pulse into a CELP coder by adding an extra codebook. The residual quality of synthetic speech was reported to be improved over the quality produced by the primitive CELP coder. Recently, the incorporation of the residual features by means of successive codebooks also gained a certain degree of

reproduction in synthesizing natural speech at 1.6 Kbit/s (Huang et al., 1990; Zhang and Chen, 1992).

Other attempts to replace the exciter by synthesized pulses appear in the work by Jambur et al. (1981) and by Chikara and Wu (1990). Among the tested pulses, the differentiated electroglottograph (DEGG) signal was found to produce good quality. Such a result occurred because the DEGG signal reflects the glottal characteristics and has a richer timbre spectrum.

In contrast to the lumping approaches, many researchers have adopted a "divide-and-conquer" strategy to design the exciter. Some divided the spectrum of the exciter into several frequency bands and examined the corresponding spectral characteristics for each band (Malkhous et al., 1979; Kwon and Goldberg, 1984; Griffin and Liu, 1988; McCree and Bjornsvik, 1981). The excitation signal was then formed by summing the subband components under the constraint that the resulting waveform must exhibit a flat spectrum. If there were only two spectrum bands to be specified, the model was often referred to as the mixed excitation since it consisted of a mixture of low-frequency pulses and high-frequency noise. If the number of divided bands matched that of pitch harmonics, this type of excitation became a superposition of sinusoids and was termed random or pseudo-periodic or either sinusoids or sinusoids (Tremecan et al., 1990). On the other hand, such a "divide-and-conquer" strategy was also considered by researchers for use in the voice domain. Sinervo (1988) parted the exciter signal into three parts, i.e., high energy pulses, a low energy smooth component, and a random noise component. Each component was acquired by using a distinctive feature. For instance, the high energy pulses were found to exhibit some similarities inherent similar to MFLP coders. After subtracting the pulses from the exciter, the smooth component was calculated by vector quantization. Likewise, the noise component was determined by codeword searching as in CELP coders. Such an approach has proven useful for speech coding in the range of 5.6 Kbit/s. Sakai et al. (1989) decomposed the exciter into a set of orthogonal functions called Zoni functions. They

claimed that the *Flow* function is superior to the *Formant* approach for modelling the vowel in the vowel space contours. However, even though both the frequency and time domain approaches offered better synthetic quality, none of the above mentioned models provided a clue to describe the glottal features parametrically.

From the above discussion, it appears that the quality of synthesized speech can be improved when we attend to the basic features of the vowel signal. Our investigation showed that the vowel was closely related to the glottal volume velocity via the glottal shaping filter (see Section 2.3). In fact, Kang and Stevens (1985) have demonstrated how to improve the quality of the pitch-started LPC vocoder through the exploitation of the amplitude and phase spectra of the vowel. It was also reported that high-quality LP synthesis could be achieved by introducing an extended filter which captured some of the glottal phase characteristics (Chang and Lee 1987; Griffin, 1988). The improvement due to the appropriate modeling of glottal source is more evident when a glottal flow model is applied to the formant synthesizers (Karnenberg, 1971; Holmes, 1973; Klatt, 1980; Pons, et al. 1989), but such parametrically important functions have not been widely considered in LP synthesis. Our primary goal is to design an efficient, accurate model to describe the vowel so that we can achieve high-quality natural-sounding speech production using such an excitation model.

2.1 Excitation Source

In a manner similar to that adopted in the traditional LP synthesizers, we classify the excitation function into two categories, i.e., voiced and unvoiced. Accordingly, two different strategies are employed to analyze and process the speech signal.

2.1.1 Voiced Segments: Excitation Policy

In Section 2.3, we have shown that the phase characteristics of a glottal flow waveform could be extracted from the residue signal. However, since the non-reflexion

level of the global flow has been lowered due to the lower sampling and integration, some models that specify the differentiated global flow are not suitable for modeling the integrated model. We therefore propose a new model to code the signal of the model. This model is described by a cubic interpolational $f(t) = \sum_{i=0}^3 a_i t^i$, which is specified within the interval [0,1] subject to three constraints listed below:

$$1. f(0) = -1. \quad (3-1)$$

$$2. f(1) = f(0). \quad (3-2)$$

$$3. \int_0^1 f(t)dt = 0. \quad (3-3)$$

where the interval has been 0, 0.5 and 1, corresponding to the global clock instant (GCI). The order of the polynomial is explicitly chosen to be six because it sufficiently describes the integrated model without causing rank deficiency.

The purpose of the constraints is as follows. The first constraint is used to normalize the magnitude of the input signal or peak. The second constraint is to ensure the circular continuity between consecutive periods. It is also equivalent to the following expression:

$$\int_0^1 f'(t)dt = 0, \quad (3-4)$$

which indicates that the $d.c.$ component in the model signal is eliminated. The third constraint is established to avoid any low-frequency modulations.

Because of these constraints, only four degrees of freedom are available in the polynomial even though seven coefficients exist. To require the polynomial coefficients under such constraints, we can introduce Lagrange multipliers and solve a set of equations as in an optimal control system. Nonetheless, the main purpose of these constraints is not

to limit the dynamics of the polynomial coefficients while carrying out the approximation. Instead, the constraints are just used to regularize the polynomial waveform. They can also be satisfied by adjusting a reference polynomial, which is calculated based on a least square fit. Here we apply a weighting function to emphasize the polynomial fitness around the GCI since this region is directly related to the primary excitation pulse. The weighting function is given by

$$W(x) = \begin{cases} 300x^2 - 60x + 3 & \text{for } 0 \leq x \leq 3 \quad \text{[I]} \\ 1 & \text{for } 3 \leq x \leq 8 \quad \text{[II]} \\ 25x^2 - 40x + 17 & \text{for } 8 \leq x \leq 1 \quad \text{[III]} \end{cases} \quad (3-5)$$

and is displayed in Figure 3-1. In practice, the weighting function can also reduce the chance of rank deficiency while we perform the polynomial fit.

Once we obtain the reference polynomial, the first constraint can be achieved by normalizing all the coefficients with respect to C_0 , i.e.,

$$C_i = -\frac{C_i}{C_0} \quad \text{for } i = 0, 1, 2, 3, 4, 5, 6 \quad (3-6)$$

The second constraint can be satisfied by setting a value r close to 1 such that $g(r) \approx 1$. Accordingly, the polynomial coefficients are revised as

$$C_i = C_0 r^i \quad \text{for } i = 1, 2, 3, 4, 5, 6 \quad (3-7)$$

The solution regarding the third constraint came out to be a procedure for removing the d.c. level. We can modify the constant C_0 to accomplish this requirement

$$C_0 = -\sum_{i=1}^6 \frac{C_i}{i+1} \quad (3-8)$$

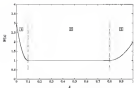


Figure 3-1 Plot of the weighting function $W(x)$

Then, the residual integral for a period becomes

$$\int_0^T f(t)dt = \sum_{n=0}^L \frac{F_n}{n+1} = 0 \quad (3-6)$$

It is important to note that when applied to the initiative polynomial the above mentioned adjustments shall be arranged as Eq. (3-7), then Eqs. (3-8) and (3-9) in order to prevent any further conflict among the constraints.

3.2.1.1 Vector quantization

Like many other glottal source models, the polynomial model only provides a rough description of the glottal phase characteristics. The lack of detailing of the glottal phase may lead to a degradation in quality of synthetic speech. However, in a study concerning the influence of glottal flow waveforms on the quality of voiced synthetic speech, Rosenburg (1977) concluded that only gross source features are required to preserve the quality whereas temporal and spectral details are less important. Assurances regarding the phase characteristics were further supported by other researchers (Arai and Dorai, 1978; Shikata, 1984). Their results lead to a speculation that the glottal waveform acquired by our model may provide sufficient discriminatory information in order to synthesize good quality speech. It is noted that vector quantization techniques have demonstrated good performance in compressing LP features with a relatively low bit rate (Landy et al., 1984; Goup, 1984). We believe that the glottal phase characteristics portrayed by our source model could form an compact yet an appropriate vector quantizer, at least in terms of perceptual quality.

In a general sense, the quantization is a process for converting a continuous-amplitude sample into one of a set of discrete-amplitude sample variables for storage and communication in a digital system. The process is known as scalar quantization, if each individual sample is quantized independently. When a block of samples, usually defined as a vector, is quantized jointly, the process is termed vector quantization.

Given a K -dimensional Euclidean space R^K , a vector quantizer considered as a partitioning, R^K into a finite subset T of R^K , where $T = \{y_i, i=1,2,\dots,N\}$ is the set of representative vectors and N the number of vectors in T . The set T is called a *codebook*, and its elements are called *codewords* or *codevectors*. In principle, the codeword y_i is chosen to minimize the average distortion for each quantized cell. The distance between any input vector and its corresponding codeword is known as the *distortion*. Once these codewords are established, any input vector is then assigned to a particular codeword based on minimum distortion for optimal representation. More specifically, pattern vector x is encoded by codeword y_i if the distance between these two vectors is less than the distance to any other codeword, i.e.,

$$d(x, y_i) < d(x, y_j), \quad j \neq i, j = 1, \dots, N \quad (3-11)$$

where the function d denotes the distance measure, and N is the number of codewords. A major advantage with the vector quantizer is that it often reduces the number of bits required to represent the input vector under a specific distortion measure. Indeed, this advantage can be formally proven through mathematical derivations. According to the Shannon rate-distortion theory, the vector quantizer always achieves higher data compression ratio than any coding scheme based on the scalar quantizer for a given transmission bit rate. Because of this, during the past decade, the vector quantization has received much attention as a data compression technique for encoding data or information intensive fields such as image and speech signals.

A vital step in establishing the vector quantizer is generation of an optimal codebook. Here the word "optimal" stands for having minimum distortion. The accomplishment of this step requires a criterion to quantify the Euclidean space and a distortion measure to define the performance of a quantizer. There are two dominant criteria commonly adopted for vector quantizers, namely, either minimizing the average quantization error or maximizing the codebook entropy defined as

$$H = - \sum_i P_i \log_2(P_i) \quad (3-11)$$

where P_i is the relative frequency with which codeword i is used to encode the sample values. While it may seem intuitive to demand a quantizer to minimize the average distortion, the most efficient way to quantize the vector space is to let each quantized cell (also known as "cluster" in some literature) consist of the same energy. Conceptually, minimizing the average quantization error can be viewed as a scheme performing a geometric division of the vector space, while maximizing the energy is a scheme to achieve a popular division of the vector space. Our philosophy of quantizing the vector space is to minimize the quantization error but at the same time to maximize the selected frequency of each codeword. It appears that the need of more relative distortion serves as a proper criterion for cluster splitting because this criterion takes both geometric and popular division properties into account (Tou and Gonzalez, 1974; Lloyd and Sederberg, 1979). Under such a criterion, clusters containing excessive training sample vectors are more likely to be split despite their inter-cluster distortions are low. Hence the codebook space is not wasted in accommodating isolated pattern vectors of global phase signals.

A perfect partition for the pattern space may be quite difficult to accomplish, although it is theoretically achievable when the distortion measure is specified and the probability density function of input vectors is known. Such a difficulty, however, can be circumvented by making use of long training sequences that approximately represent the probability density function. Thus, if the vector process is ergodic and stationary, averaging the distortion for a large amount of training vectors is equivalent to applying the probabilistic model in the underlying process. Since each vector is mapped into only one particular codeword, the codewords themselves may be established through clustering techniques. In fact, the optimal codeword is just the centroid of its associated clusters subject to a selected distortion measure. This implies that the cluster analysis algorithms or pattern recognition techniques, such as K-means, ISODATA, GMM, and some neural-net techniques can be

used to compare the existing vectors into clusters or, equivalently, to determine the hyperplane partitioning the domain (Rao and Goncalves, 1994; Rao, 1999; Rao, 1999).

3.2.1.2 Maximum descent algorithm

In this study we generate a 30-entry codebook using a maximum descent algorithm (Rao and Chao, 1991). We note that the size of the codebook is just sensitive. The number n , in general, determined by the transmission system and the desired-compression ratio.

The maximum descent algorithm states that the clusters are chosen one at a time attempting to achieve a maximum reduction of the sum of the distortions. As illustrated in Figure 3-2, we begin the splitting routine by placing all vectors in a global cluster. After forming the first two clusters, we compare the reduction functions, R_1 and R_2 , if the two are chosen and then split the one giving the larger reduction. To generalize the preceding procedure, let us consider the case of forming $n+1$ clusters based on a set of n clusters. The cluster R_n (or that it splits into two new clusters if R_n is the largest among all the R_i 's of the n clusters. Hence the set of $n+1$ clusters is the one that gives the maximum descent distortion when formed from the set of n clusters. The algorithm iterates until the desired number of clusters is obtained. Finally, the centroids of the clusters are taken as the codewords.

The advantages of using the maximum descent algorithm include: (1) computations time is significantly reduced since only the R_i 's of the two newly formed clusters need to be computed while all other clusters have been calculated in the previous iteration; and (2) empty clusters are prevented since it is as possible for a single-member cluster to be chosen for splitting.

3.2.1.3 Cluster splitting

Split (with subvector) represents the centroid of a specific cluster; the size of the codebook equals the number of clusters partitioned in the pattern space. We adopt a splitting technique to carry out the cluster partition. This technique, in general, is not guaranteed to

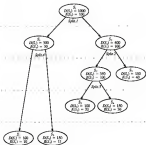


Figure 3-2. Cluster splitting using the Maximum Deviant method.
 DCL —sum of distances.
 $DCLD$ —deviations of distances due to cluster splitting.

provide an optimal solution, but it gives satisfactory results even with a binary search coding scheme (Bout et al., 1980). Steps for splitting each given cluster are summarized as follows.

Step 3: Assign the initial centroids by using the extreme-point approach, which will be discussed in Section 3.2.1.3.1.

Step 4: Partition the cluster vectors on the basis of minimum distortion, i.e.,

$$\begin{aligned} \text{if } d(x_i, c_1) < d(x_i, c_2), \quad x_i &\in S_1; \\ \text{otherwise,} \quad x_i &\in S_2. \end{aligned}$$

Step 5: Obtain the new centroids by

$$\begin{aligned} c_1' &= \frac{1}{N_1} \sum_{x_i \in S_1} x_i \\ c_2' &= \frac{1}{N_2} \sum_{x_i \in S_2} x_i \end{aligned}$$

where N_1 and N_2 are the number of vectors assigned to S_1 and S_2 , respectively. The superscript i denotes the number of iterations.

Step 6: Calculate the reduction of distortion R_i^d due to splitting as

$$R_i^d = D(X_i) - [D(X_{i-1}) + D(X_{i-1})],$$

$$\text{if } \frac{R_i^d - R_{i-1}^d}{R_i^d} > 10^{-5}$$

then Go to Step 3;

else Terminate.

The success of the cluster analysis, in general, will be affected by three factors, namely, the initial centroids, distortion measure, and the geometric properties of the training vectors. The geometric properties reflect the distribution of feature patterns and can be adjusted by properly selecting the training vectors. At this stage, we can assume that the

selected training vectors are necessary both in completeness and equilibrium. Thus, we are concerned only with the initial constraints and distance criterion.

3.2.1.2.1 Initialization of constraint

Several methods for determining the initial codebook exist. We may simply choose the first two training vectors as our initial constraints, similar to the manner used in the K-means method. However, simply choosing the first two vectors will not produce an accurate result if these two vectors are close to each other. Intuitively, one would like these two vectors to be well separated. We therefore, assign the two initial constraints using the following approach. Let $\{x_1, x_2, x_3, \dots, x_N\}$ be the N sample vectors. The mean vector is given by

$$\bar{x}_0 = \frac{1}{N} \sum_{i=1}^N x_i \quad (3-12)$$

Using \bar{x}_0 as a reference vector, we then find a vector x_{i_0} that is furthest from \bar{x}_0 . That is,

$$d(x_{i_0}, \bar{x}_0) > d(x_i, \bar{x}_0) \quad \text{for } i \neq i_0, \quad i, i_0 = 1, \dots, N \quad (3-13)$$

This vector x_{i_0} is selected as one of the constraint vectors. The other is determined by searching for the vector that is furthest from x_{i_0} .

3.2.1.2.2 Distance criterion

As mentioned earlier, the feature space consists of the polynomial coefficients. We adopted the Euclidean distance as the distance measure, which is defined as

$$d_f = \int_0^1 (g(x) - \hat{g}(x))^2 dx \quad (3-14)$$

where d_f is the resulting criterion of two arbitrary polynomials, $g(x)$ and $\hat{g}(x)$. The constraint, $\hat{g}(x)$ of a cluster, S_{i_0} is chosen as

$$J_0(x) = \frac{1}{N_0} \sum_{i \in \mathcal{N}_0} J_0(x_i) \quad (3-13)$$

where N_0 is the number of vertices inside Ω_0 . Thus, the sum of distances DU_0 for the cluster Ω_0 is given by

$$DU_0 = \sum_{i \in \mathcal{N}_0} \int_0^1 (g_i(x) - J_0(x))^2 dx. \quad (3-14)$$

Let P_i denote the vector of the polynomial coefficients, $g_i(x) = \hat{g}_i(x)$, in a descending order. The polynomial multiplication of $g_i(x) - \hat{g}_i(x)^2$ is equivalent to convolving P_i with itself, i.e., $P_{2i} = P_i * P_i$, where $[*]$ denotes the convolution operator and P_{2i} is the coefficient sequence of resulting polynomial. After solving the integral function, Eq. 3-14 becomes

$$DU_0 = \sum_{i \in \mathcal{N}_0} \sum_{n=0}^N \frac{P_{2i}(n)}{(n+1-\alpha)} \quad (3-15)$$

where n is the number of coefficient of P_{2i} . In our case, $n = 13$.

3.2.1.4 Codebook Training

In order to reflect the source variation causality factors such as stress and intonation, we use sentences instead of isolated words for training the global codebook. The selected sentences are (1) "We were away a year ago," spoken by 18 subjects, and (2) "Early one morning a man and a woman walked along a road early here," spoken by 4 subjects. In both cases, the numbers of both male and female subjects are equal. The data type is shown in Table 3-1. The resulting codebook is given in Table 3-2. The inclusion of words, as in the second instance, is intended to compensate for the deficiency of the all-pole model by introducing more (non-formant) characteristics to the source model. Although the set of training samples does not consist of all possible voiced sounds, source properties are still considered representative since the suprasegmental loading effects are removed by the stress filter and source characteristics are presumably the only remaining ingredients.

Table 3-1 Database for codebook training.

Initials	Sex	# of pitch periods	use source	Initials	Sex	# of pitch periods	use source
STB	M	400	(1)	CLN	F	290	(1)
DB	M	313	(1)	ADM	F	303	(1)
DMB	M	356	(1)	CRG	F	337	(1)
POD	M	323	(1)	FRB	F	288	(1)
DAAR	M	335	(1)	BRW	F	337	(1)
BOC	M	313	(1)	BLM	F	270	(1)
WCS	M	341	(1)	MSH	F	264	(1)
TLB	M	323	(1)	PLN	F	336	(1)
DOH	M	375	(1)	CAP	F	361	(1)
MJB	M	337	(1)	LAD	F	361	(1)

3.2.2 Unvoiced/Silence Segments: White Noise

For simplicity, we treat silence as unvoiced speech, since the power level of the silence segments is so low that any modeled errors can be attributed to background noise. Similar to the idea adopted in voiced excitation, a stochastic codebook is used as the codebook source for unvoiced speech. This implies that the residue is simulated using a finite number of excitation sequences selected as a given safety-interval. The use of such excitation sequences is motivated by the CELP codes, of which the stochastic codebook has been known to produce better unvoiced speech than voiced speech for low bit coding (Schalinski and Lacroix, 1989). Just as contrast to the fundamental structure of the CELP codes, the randomly used long-term predictor is dropped here since pitch harmonics are unnecessary in unvoiced speech.

Basically, the size of the codebook is determined by three factors, namely, the transmission rate, computational complexity and frame update rate. Due to the lack of an

Table 3-2. Content of global codebook.

Codevector	C_0	C_1	C_2	C_3	C_4	C_5	$C_6 \times 10^3$
1	-1.832	2.936	-1.932	2.254	-0.939	0.040	-0.000
2	-0.753	2.413	-0.007	1.713	-0.190	0.015	-0.000
3	-0.000	0.369	-0.793	0.000	-0.049	0.040	-0.000
4	-0.000	0.713	-0.700	0.410	-0.120	0.034	-0.000
5	-0.409	1.076	-1.000	1.004	-0.203	0.038	-0.000
6	-0.000	0.000	-0.170	0.000	-0.180	0.020	-0.000
7	-0.000	2.180	-0.170	1.494	-0.210	0.045	-0.000
8	-0.007	1.200	-1.004	0.000	-0.240	0.070	-0.000
9	-0.004	0.000	-1.000	0.000	-0.040	0.010	-0.000
10	-0.000	1.000	-1.000	0.000	-0.170	0.030	-0.000
11	-0.700	0.000	-1.000	0.000	-0.180	0.010	-0.000
12	-0.000	0.000	-0.710	0.000	-0.100	0.030	-0.000
13	-0.000	1.000	-1.000	1.000	-0.200	0.000	-0.000
14	-0.000	2.000	-1.000	1.000	-0.110	0.030	-0.000
15	-0.000	1.000	-0.410	1.000	-0.010	0.000	-0.000
16	-0.000	0.000	-0.000	1.000	-0.200	0.040	-0.000
17	-0.000	0.000	-0.000	0.000	-0.200	0.040	-0.000
18	-0.000	1.000	-1.000	0.000	-0.200	0.000	-0.000
19	-0.000	1.000	-1.000	0.000	-0.200	0.000	-0.000
20	-0.000	1.000	-1.000	1.000	-0.200	0.000	-0.000
21	-0.000	0.000	-0.000	0.000	-0.000	0.010	-0.000
22	-0.000	1.000	-0.000	1.000	-0.200	0.000	-0.000
23	-0.000	1.000	-1.000	0.000	-0.000	0.000	-0.000
24	-0.000	1.000	-0.000	1.000	-0.200	0.000	-0.000
25	-0.000	0.000	-0.000	0.000	-0.100	0.020	-0.000
26	-0.000	0.000	-1.000	0.000	-0.000	0.000	-0.000
27	-0.000	0.000	-1.000	0.000	-0.000	0.000	-0.000
28	-0.000	1.000	-0.000	1.000	-0.000	0.000	-0.000
29	-0.000	0.000	-1.000	0.000	-0.000	0.000	-0.000
30	-0.000	0.000	-0.000	0.000	-0.000	0.000	-0.000
31	-0.000	2.000	-0.000	1.000	-0.000	0.000	-0.000
32	-0.000	0.000	-1.000	0.000	-0.000	0.000	-0.000

Note: C_i denotes the i th coefficient of the polynomial.

appropriate criterion for characterizing performance, we synthetically create the stimuli for a limited duration (30 samples or 150 Hz) by the use of 256 codevectors. The type of codebook population is not a crucial factor from a perceptual point of view: experiments with Gaussian, sparse and binary-value $[-1, 0, +1]$ codebooks have been reported to produce similar synthesis quality (Troxen et al., 1990). However, since the probability density function of the synthesized residual is nearly Gaussian, we will employ a Gaussian noise generator to establish the codevectors. For each codevector, specializations of its content are only for the purpose of reducing the computational effort, which is necessitated by the shaping process in codevector searching (Kang et al., 1990; Gahleitner et al., 1992). This kind of computational burden can also be alleviated by other means such as the complex value decomposition (PVD), frequency domain and autocorrelation approaches (Troxen and Aul, 1992). In our experiments, the autocorrelation approach is adopted to facilitate the comparison. Some relevant details will be given in Chapter 4.

As mentioned in the previous section, samples for each codevector are drawn from a Gaussian noise generator, but we employ three schemes to establish the codebook.

Scheme 1 (64-codes) — Each codevector contains 16 non-zero samples.

The positions of non-zero samples exhibit a uniform distribution from 1 to 56.

Scheme 2 (64-codes) — The positions are the same as Group 1 except that 32 out of 56 samples are non-zero.

Scheme 3 (128-codes) — Every sample is taken from a Gaussian noise generator.

The sparse codevectors in Schemes 1 and 2 are used to enhance the spiky nature of the residual so that the stochastic codebook can also be applied to synthesize the natural sounds as well as phonemes. This concept is very similar to that proposed by Kang and Brown (1981), who introduce a few sparse spikes into the synthesized excitation in order to obtain satisfactory pleasant speech.

CHAPTER 4 SPEECH ANALYSIS/THESIS EVALUATION

In Chapter 2, we discussed in detail how to interpret the acoustic features of speech signals within the linear source-filter theory. The power of LP techniques for performing the feature extraction suggests that a high-quality LP synthesizer could be achieved if these features were appropriately modeled and accurately estimated. Hence, in Chapter 3 we discussed source modeling. The results, known as the *linear source excitation*, was analyzed either by glottal impulses for voiced speech or excitation sequences for unvoiced speech. Both types of excitations were further formulated into two specific codebooks. The reader can advantage Chapter 2 as an *acoustical study* of the speech signals and Chapter 3 as an *examination* of the glottal source. The information obtained from these chapters can now serve as a *driving* a synthesis model capable of producing high-quality natural-sounding speech.

In this chapter, we present a new model which includes many improved features such as the manipulation of LP coefficients, turbulent noise and source noise estimation. The parameters in this model are obtained by the analysis-by-synthesis procedure, in which the analysis denotes the process of estimating the parameters that characterize the speech signal and the synthesis denotes the process of replicating the speech signal by controlling and updating those parameters under the supervision of the speech production model. We will describe our methods and strategies in dealing with these issues. While the performance of this model is evaluated by judging its ability to produce natural speech, we also discuss the results of informed listening tests.

4.1 Analysis Scheme

The speech production model employed in this study is depicted in Figure 4-4. Except for the excitation source, the model mimics the basic structure of the pitch-excited LP synthesizer. In addition to an all-pole filter, other parameters required by the model comprise a voicing decision, normalized vocal gains, codeword values and Global Closure Events (GCEs) for the voiced speech.

In general, a pitch synchronous approach is preferred for speech processing not only because it provides better human engineering (Kendall-Smith and Childers, 1988) but also because it facilitates the synthesis work. To implement such an approach, we need to locate every GCE accurate before computing the LP coefficients. The difficulties of identifying GCEs complicate the feasibility and reliability of the implementation, making the pitch synchronous analysis practically unattractive. Thus, we decide to use a frame-based method to compute the LP coefficients. In our case, the speech synthesis pitch synchronously after determining the pitch period.

Since the speech signal is sampled at 1000/s, a linear predictor of 13th order is chosen to account for the spectral characteristics of the glottal source (3 poles) and vocal tract (10 poles). The filter coefficients along with the residue are derived successively using an orthogonal covariance method (Wang and Whiting, 1992), performed once per frame sequentially throughout the input speech. The frame size is 20 sfs with an overlap of 5 ms between any two consecutive frames. For each frame, the LP gains are estimated by adjusting the power of the residue so that of the speech signal. The residue in the overlapped area is obtained by weighting the forward and backward overlapping sequences with decreasing and increasing trapezoidal windows respectively and adding them together.

$$e(i) = \frac{R}{N+1} \frac{1-i}{1} r_1(i) + \frac{R}{N+1} \frac{1}{1} r_1(i), \quad i = 1, 2, 3, \dots, N \quad (4-1)$$

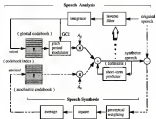


Figure 4-1. Proposed speech production model.

where $q(k)$, $x(k)$ denotes the forward and backward reaction signals respectively, $x(k)$ is the resulting reaction signal for the overlapped sets of length M .

4.1.1 Orthogonal Covariance Method

Consider a digital signal with the following sequence, $\{x_1, x_2, \dots, x_p, \dots, x_{m+1}\}$. The linear prediction of the current sample is described as a linearly weighted summation of past samples, i.e.,

$$\hat{x}_n = \sum_{k=1}^m a_k x_{n-k} + e_n \quad (4-2)$$

where the a 's are the coefficients of the LP predictor with order m , and the e 's are the prediction errors. Expressing the equations above in a matrix form, we have

$$\begin{bmatrix} x_1 & x_2 & \dots & \dots & \dots & x_m \\ x_2 & x_3 & & & & x_{m+1} \\ x_3 & x_4 & & & & \vdots \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ x_m & x_{m+1} & \dots & \dots & \dots & x_{m+m} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ \vdots \\ a_m \end{bmatrix} = \begin{bmatrix} x_{m+1} \\ x_{m+2} \\ \vdots \\ \vdots \\ x_{m+m} \end{bmatrix} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ \vdots \\ e_m \end{bmatrix} \quad (4-3)$$

For the convenience of illustration, vector notation is employed in the following derivations. We define the \mathbf{X}_n as the n th column vector of the matrix \mathbf{X} , \mathbf{A} as the vector of the LP coefficients, and \mathbf{E} as the vector of prediction error. Thus, Eq. (4-3) becomes

$$[\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots, \mathbf{X}_m] \mathbf{A} = \mathbf{X}_{m+1} - \mathbf{E} \quad (4-4)$$

By assuming that the prediction error is negligible, we drop the term \mathbf{E} and determine \mathbf{A} by multiplying the pseudo inverse of \mathbf{X} on both sides of Eq. (4-3). It may be shown that the obtained result is the same as that derived by a covariance method, because the error is minimized over a specified interval.

It can also be shown that a certain degree of efficiency could be gained by reformulating the foregoing computation as follows. Suppose we now decompose the matrix \mathbf{A}_{p+1} into p orthogonal vectors \mathbf{V}_i^T 's using the Gram-Schmidt method. The set of the orthogonal vectors is

$$\mathbf{V}_{p+1} = \mathbf{A}_{p+1} - \sum_{i=1}^p c_i^T \mathbf{V}_i \quad (4-3)$$

$$\text{where} \quad c_i^T = \frac{\mathbf{A}_{p+1}^T \mathbf{V}_i}{\|\mathbf{V}_i\|^2} \quad (4-4)$$

and the superscript T denotes the transpose operator. Arranging the orthogonal expansion in a matrix form, we obtain

$$\begin{bmatrix} 1 & 0 & & & 0 & 0 & 0 & 0 \\ a_{11}^T & 1 & & & 0 & 0 & 0 & 0 \\ a_{21}^T & a_{22}^T & 1 & & 0 & 0 & 0 & 0 \\ & & & \ddots & & & & \\ a_{p1}^T & a_{p2}^T & & & 0 & 1 & 0 & 0 \\ & & & & & & \ddots & \\ a_{p+1,1}^T & a_{p+1,2}^T & & & & & & 1 \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \\ \mathbf{V}_3^T \\ \vdots \\ \mathbf{V}_p^T \\ \mathbf{V}_{p+1}^T \end{bmatrix} = \begin{bmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \\ \mathbf{V}_3^T \\ \vdots \\ \mathbf{V}_p^T \\ \mathbf{V}_{p+1}^T \end{bmatrix} \quad (4-5)$$

Through several algebraic manipulations, the row vector \mathbf{V}_{p+1}^T on the right-hand side of Eq. (4-5) can be shown as

$$\mathbf{V}_{p+1}^T = \sum_{i=1}^p c_{p+1,i}^T \mathbf{V}_i^T + \mathbf{V}_{p+1}^T = \mathbf{C}_{p+1}^T (\mathbf{C} \mathbf{V}_{p+1}) \begin{bmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \\ \mathbf{V}_3^T \\ \vdots \\ \mathbf{V}_p^T \end{bmatrix} = \mathbf{V}_{p+1}^T \quad (4-6)$$

where \mathbf{C} is the matrix containing the c_i^T 's on the right-hand side of Eq. (4-5); $(\mathbf{C} \mathbf{V}_{p+1})$ is a upper-left submatrix of the matrix \mathbf{C} with a rank p ; and \mathbf{C}_{p+1}^T is the $(p+1)$ th row vector of the matrix \mathbf{C} with only the first p -elements included. Compared to Eq. (4-3), it is found

intuitively that the coefficient vector A is equivalent to the term

$$A' = C_{p+1}^{-1} / (C_{p+1}^{-1} C_{p+1}^{-1}) \quad (4-8)$$

and the vector \hat{Y}_{p+1} is just the estimated error (or the residual) of the p th order LP filter

A major advantage of using the orthogonal expansion is that the matrix inverse of C can be achieved by a back-substitution procedure (Ning and Whiting, 1992). Another important advantage emerges from the error vector \hat{Y}_{p+1} , which is known as a principal component to determine the order in many methods. In speech processing, the importance of selecting a correct order can be explained in terms of formant characteristics. A filter with a lower order tends to disregard insignificant formants or to merge two adjacent ones, whereas a higher order filter runs the possibility of producing spurious formants. The resulting incorrect formants may lead to perceptible errors in both cases. Thus, a variable-order predictor is always preferred, for it adapts the spectral variation of missing speech. Apart from this reason, such a model can also reduce the unnecessary bandwidth of lower orders are frequently chosen.

There are two widely accepted order measures, namely, the Akaike information criterion(AIC) (Akaike, 1974) and minimum description length(MDL) criterion (Schwarz, 1978) which can be obtained prior to estimating the LP coefficients. Two other methods proposed in the recent literature are the Predictive Least Squares (PLS) (Wan, 1992), and the Iterative Algorithm of Singular Value Decomposition (IASVD) (Kammarathani and Yoo, 1992). A comparison of the performance of the four methods indicated that the IASVD had the highest success rate in order selection followed by MDL, AIC and PLS (Kammarathani, 1992). If we take into account computational efficiency, the MDL term can be a proper choice to work with the orthogonal covariance method. Eventually the selected order is the one that minimizes the MDL function given by

$$REDC(i) = N \ln \ln \left(\frac{P_i^2}{P_1} \right) + i \ln(N) \quad (4-10)$$

After we determine the optimal order p , the LP-coefficients a are derived from Eq. (4-9).

4.1.1.3 Vowel Classification

Because there are two types of excitation functions in the proposed model, the first step toward speech analysis is a vowel decision. The basic principle of our method is rather simple. If the energy of the underlying signal is below a specified value, the signal is classified as silence (Changbali and Thomas, 1986; Chikara et al., 1989a). Otherwise, we examine its spectral slope by calculating the first-reflection coefficient. The signal is classified as voiced speech if the first reflection coefficient is larger than 0.3. Unvoiced speech is the result when the previous two methods failed.

Unlike other algorithms, the correct rate of classification is not strictly required because the incorrect decisions does not lead to serious perceptual errors. For example, the cross classification between unvoiced and silence is not critical since both share the same excitation functions and the quantization error in the silence was always be ignored. Also, the speech signal with a random spectral tilt (e.g., the first-reflection coefficient is around 0.3) often exhibits mixed characteristics of both types of excitations. Therefore, either voiced or unvoiced classification is acceptable for synthesizing such a speech signal.

4.1.1.4 Identification of Global Closure Interval (GCI)

A reliable identification of the GCI is essential for colored-word searching and speech synthesis since both are performed on a peak-pick-by-peak period basis. Procedures of the GCI classification algorithm can be summarized in two steps: (1) peak extraction, and (2) peak picking. That is, we determine the location of global closure after extracting the peak period.

It has long been noted that the sharp peaks in the random signal generally coincide with the GCIs for a wide variety of voiced sounds. Choosing the largest peak of the random

signal for many voiced words is a useful method for determining the GCB (Niel and Blamont, 1971; Ananthapadmanabha and Nageswaraswara, 1979). For words that are not rich in harmonic structure or that lack distinctive glottal closure may fail to have large peaks in the residual. Furthermore, the true peaks may be obscured by other spurious peaks due to background noise and modelling errors. Perhaps the easiest way to circumvent this drawback is to apply a lowpass filter to reduce the influence of the spurious peaks. However, this results in smoothing and thereby decreases the sharpness of the real peaks, so we can only diminish the influence of noisy components to such an extent that the true peaks are not masked out. To avoid the phase shift of the peaks, we perform the lowpass process by a zero-phase filter, i.e., by first passing the residual signal forward then running it back through the same filter (Oppenheim and Wilsky, 1983). The Z -domain representation of the employed filter in our experiments is chosen to be

$$H(z) = \frac{1}{(1 - 0.9z^{-1})(1 - 0.9z^{-1})^*} \quad (6-1)$$

Once the residual signal is lowpass filtered, a segment of 312 samples centered at the current frame is extracted using a framing window. This windowed segment $x(n)$ is then transformed to a sequence $P_1(n)$ similar to the expression by

$$P_1(n) = |FFT(FFT^{-1}(x(n)))| \quad (6-2)$$

where FFT and FFT⁻¹ stand for the fast Fourier transformation and its inverse operation, respectively, and $| \cdot |$ denotes the magnitude.

Like the pitch estimation procedure outlined in the original method, we choose n as the pitch period if

$$P_1(n) \geq P_1(m) \quad m = 25, 50, \dots, 125. \quad (6-3)$$

The value m could be a multiple of the real pitch period, however. In our program, a sample

check is given as follows. We first look for the position of the largest value within the range $[25, m-25]$, i.e.,

$$P_1(l) \geq P_1(k) \quad \text{for } l \neq k \text{ at } l, k = 25, \dots, m-25. \quad (3-4)$$

If the following condition exists

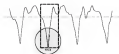
$$P_1(l) \geq 3P_1(m), \quad (3-5)$$

then l is adopted as the new pitch period. Otherwise, the pitch period is assumed as m .

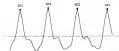
After finding the pitch period, we begin with the search for the largest negative peak in the resampled signal. Due to the fact the peak has been smoothed by the zero-phase lowpass filter, we enhance the accuracy of peak picking by approximating the curve on both sides of the negative peak by two straight lines ranging from the peak value to one-third of this value (highlighted by the dotted line in Figure 4-3(a)). The intersection of the two lines is chosen to be the first GCI. A small interval of samples (~ 4 bins) around the first GCI is used as a template (as shown as a dashed box) to discriminate other peaks within the same frame. Peaks located before or after this GCI with approximately one pitch period are examined by computing the correlation between the template and the waveforms around the peaks. Positions that lead to largest correlation are then selected as other GCIs. This procedure continues until the searching range is out of the current frame by 50 samples.

The overall computation above costs 2 FFT's and several comparisons. An economical approach for performing the whole process is to decimate the signal $x(n)$ by a factor of 2 and then to perform an interpolation on $P_1(n)$ to maintain such a decimation. Because of the lowpass filtering, the lowpass decimation can be carried out by choosing every other samples of $P_1(n)$ without causing aliasing.

A complete example of the GCI identification is illustrated in Figure 4-3(b). In this example, it appears that the GCIs can be directly identified by picking the negative peaks of the lowpass filtered waveform. But, as we mentioned earlier, the peaks on the original signal



(a)



(b)

Figure 4-2 Illustration of OCT identification. (a) lowpass filtered radar signal, (b) cross-correlation for the radar signal with a template stored in the database (see 3d)

are not always sharp and distinctive. The inherent weak mode signal in the measured residues plays a role in helping to reduce potential errors. Similarly, the acquisition of GCT controls the control frame is instantaneous but necessary, because the state information can be used to prevent erroneous GCTs in future boundaries.

4.1.3 Code-word searching

Depending on the coding conditions, there are two different codebooks prepared to compress the systematic streams. Although the basic idea of code-word searching for these two codebooks is the same, i.e., selecting a optimum code-word that achieves a minimum error subjective distance metric, the individual implementations are somewhat different due to their intrinsic characteristics.

4.1.3.1 Fixed-resolution, global codebook

The searching process for the optimal global code-word requires that the integrated residue and the polynomial waveform are of the same length. We assume the maximum allowable length for one pitch period to be 256 ms. Thus, if we encode every polynomial waveform with such a maximum length, then the integrated residue of one pitch period can always be interpolated to the maximum length using the FFT method. Taking advantage of the symmetry of the Fourier transformation, we compute the correlation coefficient, $\eta_{ab}(l)$, between the a th polynomial waveform, $p_a(x)$, and the integrated residue of the a th period, $d_a(x)$, by

$$\eta_{ab}(l) = \frac{\text{real} \left[\sum_{n=0}^{\frac{L}{2}-1} D_a(n) G_b^*(n) \right]}{\left[\sum_{n=0}^{\frac{L}{2}-1} D_a(n) D_a^*(n) \right]^{1/2} \left[\sum_{n=0}^{\frac{L}{2}-1} G_b(n) G_b^*(n) \right]^{1/2}} \quad (4-6)$$

$$= \frac{\text{real} \left[\sum_{j=0}^{\lfloor N/2 \rfloor} \tilde{D}_m(j) \tilde{D}_m^*(j) \right]}{\left[\sum_{j=0}^{\lfloor N/2 \rfloor} \tilde{D}_m(j) \tilde{D}_m^*(j) \right] \left[\sum_{j=0}^{\lfloor N/2 \rfloor} \tilde{D}_m(j) \tilde{D}_m^*(j) \right]} \quad (4-6)$$

where $\tilde{D}_m(j)$ is the FFT sequence of the interpolated $d_m(t)$, $\tilde{D}_m^*(j)$ is the FFT sequence of $\tilde{d}_m^*(t)$, $\lfloor \cdot \rfloor$ denotes the ceiling function, and $(\cdot)^*$ denotes the complex conjugate. It is noted that the mean values of $\tilde{d}_m(t)$ and $\tilde{d}_m^*(t)$ are zero and, therefore, play no role in computing the correlation coefficients. We reify this consequence by dropping the DC term during the multiplication of two FFT sequences. The second equality in Eq. (4-6) is due to the fact that the interpolated FFT sequence of $\tilde{D}_m(j)$ is zero when $k \neq \lfloor N/2 \rfloor$. The spectral weighting filter, which is commonly used in the CRLF codes, also can participate in the equation above. This is because our distance measure is applied to the integrated window, which emphasizes only on the global phase characteristics at the low frequency region.

Since the global waveform varies relatively slowly compared to the changes of the rapid inter-frame features, one codeword is used to describe the global variation for each coded frame. We further define the cross-frame similarity function, $H(k)$, as the sum of $q_m(k)$ along one frame,

$$H(k) = \sum_{m=1}^M q_m(k) \quad (4-7)$$

where M is the total number of the pitch periods in one frame. The codeword that leads to the maximum correlation similarity is chosen as the representative for the coded frame.

4.1.3.3 Hierarchical exhaustive searchcodebook

In this category, the codeword searching method is that commonly used in CELP coders. The remainder of this section provides a brief discussion of the CELP algorithm and the implementation method that achieves a fast codeword searching.

4.1.3.3.1 CELP algorithm

The CELP algorithm was first proposed by Schroeder and Atal in 1984. A rapid development revolutionized this algorithm in the last 1980's. The CELP coder represents a breakthrough in speech coding for a narrowband speech signal at a rate as low as 8 kbps but still produces a satisfactory quality. The basic concept for the class of CELP coders can be viewed as a vector-quantization technique, which passes a finite set of candidate vectors through an all-pole predictive filter and then selects the one giving a best match response to a specific excitation waveform. However, the research that led to the development of CELP coders seemed to follow another path emerging from the MPELP coder, in which the excitation consists of a few pulses per frame regardless of whether the speech is voiced or unvoiced. The locations and amplitudes of these pulses are determined by minimizing a subjective error between the original and synthetic speech signals. The relationship between the CELP and MPELP coders can be understood by considering the multipulse excitation as a deterministic codebook consisting of excitation sequences (or codewords), each consisting of M single impulses with a different delay. Hence, searching for an optimum pulse location across the analysis frame is equivalent to searching through a set of excitation.

In the previous CELP coder (Figure 4-1), the speech signal, $x(n)$, is sampled in blocks of N samples. For each block, the synthetic speech signal is derived by filtering every excitation sequence stored in a codebook into two recursive filters (long term and short term) with a proper scaling factor. An error signal is then formed by comparing the synthetic speech with original one. Through an exhaustive search over the entire codebook,

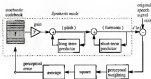


Figure 6-3. Block diagram of the CELP vocoder

the innovation sequence (along with an appropriate scaling factor) that produces the minimum mean-squared subjective error is selected to reconstruct the synthetic excitation.

The short-term predictor in the CELP coder is the well-known LP filter. The long-term predictor is an echo stage used to enhance the periodicity of the synthetic speech by replicating the stationary across consecutive pitch periods, and has been applied in open-loop and closed-loop form. In the former case the long-term predictor is directly derived from the excitation obtained by inverse filtering the original speech, while in the latter case the optimal long-term predictor is computed based on an analysis-by-synthesis procedure. Although the analysis-by-synthesis procedure does not provide much improvement of speech quality over the open-loop procedure, it opens the concept of the "adaptive codebook" or "self-excited" model, in which the codebook entry is selected as the optimization of a matching window to the recent past excitation. More precisely, each entryword is a shifted version of the previous one with one new sample changed at the end. The conceptual structure of the adaptive codebook is illustrated in Figure 4-4. As seen in the figure, the function of the pitch predictor is replaced by the adaptive codebook. Owing to the dependency of the neighboring codewords, together with a related error criterion that provides an even weighting to the codewords, fast algorithms have been exploited to reduce the inherently high-computational complexity of closed-loop procedure.

Following the formalism given by Tiedeman and Atal (1980), we now use the matrix notation as well as vector notation to illustrate the analysis-by-synthesis procedure for codeword searching. Given a codebook of L responses $C_k^{(P)}$ ($k=1, 2, \dots, L$) each of length K , the filtering operation for an innovation sequence by the long- and short-term filters can be carried out by convolving the innovation sequence with the combined impulse responses of these two filters. Written in matrix form, the filter output for the k th codeword can be expressed by

$$y^{(k)} = y^{(P)} H_k^{(P)} \quad (4-8)$$

where $y^{(k)}$ is the coding factor for the k th subword, H is an $M \times M$ matrix with the elements in the m th row and the n th column given by the two-way sample of the autoimpulse response of the filter, and $x^{(k)}$ is a M -dimensional vector with its m th component given by $x_m^{(k)}$. Since $h_{00} = 0$ for $m=0$, the matrix H can be shown as

$$H = \begin{bmatrix} h_{00} & 0 & \cdots & \cdots & 0 \\ h_{10} & h_{00} & \cdots & \cdots & \\ \vdots & h_{10} & \ddots & \ddots & \\ h_{M-1,0} & h_{M-1,1} & \cdots & h_{10} & h_{00} \end{bmatrix} \quad (9-8)$$

Let us define x to be the desired signal with its m th component given by x_m , of which the memory contribution carried over previous frames has been removed since the filter memory plays no role in the search procedure. The total squared error $J^{(k)}$, representing the difference between the desired vector x and the vector $y^{(k)}$, is defined as

$$J^{(k)} = \|x - y^{(k)} H_c^{(k)}\|^2, \quad (9-9)$$

where $\|\cdot\|^2$ indicates the squared norm of the underlying vector. The optimum code factor $y^{(k)}$ that minimizes $J^{(k)}$ is determined by setting $\partial J^{(k)} / \partial y^{(k)} = 0$, yielding

$$y^{(k)} = \left[\frac{x H_c^{(k)T}}{H_c^{(k)T} H_c^{(k)}} \right]^T, \quad (9-10)$$

and the error becomes

$$J^{(k)} = x^T x - \frac{(x H_c^{(k)T})^2}{H_c^{(k)T} H_c^{(k)}} = \|x\|^2 - \left[\frac{x H_c^{(k)T}}{H_c^{(k)}} \right]^2 \quad (9-11)$$

The best subword is obtained by selecting the index i in a exhaustive search for which the error $J^{(i)}$ is minimum (i.e., equivalently, the second term on the right hand side of Eq. (9-11) is maximum).

In principle, the error derived above spans over the entire spectrum of the synthetic speech. Due to auditory masking, the error in the high energy regions is masked by the speech signal, suggesting that the error should be reconstructed in the downward regions to reduce perceptual distortion. This idea can be easily accomplished by the use of a weighting filter $W(f)$ that attenuates the frequencies where the error is perceptually less important, and amplifies those frequencies where the error is perceptually more important.

$$W(f) = \frac{1 - \sum_{k=1}^L a_k f_k^2 z^{-k}}{1 - \sum_{k=1}^L a_k f_k^2 z^{-k}} \quad (4-11)$$

where if $k a_k f_k^2 < f_k^2$ is 1, using a_k as the LP coefficients. If f_k is set to be very close to the range $0.4-0.6$ kHz, great similar subjective results in informal listening tests (Rao and Sencoff, 1980).

Referring to Eq. (4-11), the computation required in the codeword searching contains only two terms, namely, a cross-correlation term between vectors $r^{(k)}$ and $r^{(k)}$, and an energy term corresponding to the filtered output $Wz^{(k)}$ of each codeword. The energy term is computationally simplified if the matrix multiplication $Wz^{(k)}$ is directly performed. Fortunately, many methods have been proposed for reducing the time-consuming matrix multiplication. We will now discuss the autocorrelation method, which is very efficient for fully populated autocorrelation, therefore, as illustrated above for the autocorrelation method.

4.1.3.3. Autocorrelation method

Let us first consider the energy term in the second part on the right side of Eq. (4-11). Recall that we already dropped the long-term prediction in our model. The h_k only represents the response of the impulse response of the short-term predictor filter. We rewrite the energy term in scalar notation,

$$\|M^{(R)}\|^2 = \sum_{n=1}^R \left[\sum_{m=1}^R h_{n-m} e^{j\theta_m} \right]^2 \quad (4-14)$$

Making use of the fact that the sum of the squares of the convolution of two sequences equals the cross-correlation of the autocorrelations of these two sequences, Eq. (4-14) can be simplified as

$$\begin{aligned} \|M^{(R)}\|^2 &= \sum_{n=1}^R \sum_{m=1}^R e^{j\theta_m} e^{-j\theta_n} \sum_{k=1}^R h_{n-k} h_{m-k} \\ &= R_d(M^{(R)}(0)) + 2 \sum_{m=1}^{R-1} R_d(M^{(R)}(m)) \end{aligned} \quad (4-15)$$

where

$$R_d(0) = \sum_{n=0}^{R-1-0} h_n h_{n+0} \quad (4-16)$$

and

$$R_d^{(R)}(m) = \sum_{n=0}^{R-1-|m|} e^{j\theta_{n+m}} e^{-j\theta_n} h_{n+m} h_n \quad (4-17)$$

For Eq. (4-15) to be valid, the correlation must be truncated, requiring that the impulse response of the synthesis filter is effectively zero beyond the R sample. In most circumstances this requirement will be satisfied after imposing the spectral weighting filter. If we further define the cross-correlation between h_n and h_m by

$$R_d(m) = \sum_{n=0}^{R-1-|m|} h_n h_{n+m} \quad (4-18)$$

Eqs. (4-17) and (4-15) are transformed to

$$r^{(k)} = \frac{\left[\sum_{i=1}^{N_{\text{form}}} A_i \Delta x_i^{(k)} \right]}{A_0 \Omega \Delta x_0^{(k)} \Omega + 2 \sum_{i=1}^{N_{\text{form}}} A_i \Omega \Delta x_i^{(k)} \Omega} \quad (4-25)$$

and

$$R^{(k)} = \left[r^{(k)} \right]^2 = \frac{\left[\sum_{i=1}^{N_{\text{form}}} A_i \Delta x_i^{(k)} \right]^2}{A_0 \Omega \Delta x_0^{(k)} \Omega + 2 \sum_{i=1}^{N_{\text{form}}} A_i \Omega \Delta x_i^{(k)} \Omega} \quad (4-26)$$

respectively.

From above derivation, it is easily seen that this method consumes substantial storage in computer time. The energy term can now be computed with just N multiplications per codeword. However, the price we have to pay is the storage of an additional codeword with the autocorrelation coefficients of the original codeword.

4.2 Synthesis Scheme

Speech synthesis is the procedure of reconstructing speech signals by controlling and updating the parameters of a speech production model estimated at speech analysis. The synthesis of unvoiced speech is straightforward and can be easily accomplished by exciting the time-varying all-pole filter with the gain adjusted excitation sequence sequentially. On the other hand, the synthesis of voiced speech is rather complicated because we have to suppress the syntenic excitation from the gross features of glottal pulses. Therefore, most of this section is focused on the synthesis scheme for voiced speech.

Despite many control parameters for voiced speech were estimated on a frame-by-frame basis, the corresponding synthesis can still be carried out quite systematically provided that the control parameters are properly interpolated for each pitch

period. In this section, we start from the discussion of the correspondence with respect to the glottal pulse and LP coefficients. Then, we present a method for eliminating the spectral ripples of the glottal pulse. Effects of vocal tract and room-tract interaction are discussed subsequently. Finally, a complete procedure for generating a glottal impulse is given.

3.2.1 Elimination of Glottal Phase

As mentioned earlier, only one codeword is employed to represent the glottal phase characteristics for each frame. Although large discrepancies may occur between any two adjacent pitch periods in one frame, a progressive shortness of the glottal pulse will surely occur since such discrepancies are already reflected in the codewords of different frames. This, however, results in the discontinuity of the glottal phase characteristics at the frame boundaries. Since the glottal pulse is modelled as a sixth-order polynomial, we therefore apply a lowpass filter to eliminate the rapid changes of the polynomial coefficients as follows:

(1) ER filter:

$$\hat{P}^i(k) = (1 - \alpha)P^{i-1}(k) + \alpha P_{\text{er}}(k) \quad (4-11)$$

(2) FR filter:

$$\begin{cases} \hat{P}^i(k) = \alpha P_{\text{er}}(k) + \beta P_{\text{c}}(k) + \gamma P_{\text{pr}}(k) \\ \alpha + \beta + \gamma = 1 \end{cases} \quad (4-12)$$

where $\hat{P}^i(k)$ is the polynomial for the i th pitch period, and $P_{\text{er}}(k)$, $P_{\text{c}}(k)$ and $P_{\text{pr}}(k)$ are the polynomials for the previous, current and next frames, respectively. In our program, an ER filter with the value of $\alpha=0.5$ is used since it works well in our experiments. We must note the resulting polynomial has to satisfy the three constraints specified in Section 3.2.1. Therefore, the condition where the sum of the coefficients on the right hand sides of Eqs. (4-11) and (4-12) are written as to comply with the first constraint. However, no extra consideration is necessary for the other two constraints.

4.2.2 Interpolation of LP Coefficients

As in the case of the global phase characteristics, the LP coefficients extracted from the frame-based method may exhibit undesirable discontinuities at frame boundaries. A simplification will be to linearly interpolate the LP coefficients. However, synthesis speech produced by this method may sound too smooth for speech segments with a rapid spectral transition. The plots are typical examples that suffer the drawback of the linear interpolation. This suggests that the interpolation of the LP coefficients should be "piece-wise continuous." Thus, we adopt a quadratic weighting function w_i to interpolate the LP coefficients:

$$w_i = \frac{\frac{1}{(n_p - i/2 - 1/2 + n_s - i/2 - 1/2)^2}}{\sum_{j=-1}^1 \frac{1}{(n_p - i/2 - j/2 + n_s - i/2 - j/2)^2}}, \quad i = -1, 0, 1 \quad (4-23)$$

where $|i|$ denotes the absolute value, n_p and n_s denotes the positions of the beginning and ending points of the current pitch period, and j is the number of samples in each frame. The vector of the interpolated LP coefficients, A_{interp} , is obtained by

$$A_{\text{interp}} = A_{p-1}w_{-1} + A_0w_0 + A_1w_1 \quad (4-24)$$

where A_{p-1} , A_0 , A_1 are the LP coefficient vectors of the previous, current and next frames, respectively. Figure 4-5 illustrates the linear and quadratic interpolations for an arbitrary coefficient. What we mean by "impulse linear" is clearly delineated by the quadratically interpolated curve.

One of the disadvantages of such an LP interpolation is that it occasionally moves poles outside the unit circle, implying that we have to reflect these outside poles into the unit circle in order to stabilize the synthesis filter. However, we do not consider this to be a serious problem, since the interpolation can also be done using reflection coefficients, autoregressive functions, cross-sectional areas, for all of which the stability criterion is

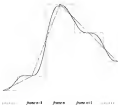


Figure 4-5. Interpolation with respect to one of the LP coefficients along several frames. (dotted line: without interpolation; solid line: with linear interpolation; dashed line: with proposed quadratic interpolation.)

satisfied. Moreover, the proposed LP interpolation conflicts with the use of a variable-order filter. Apparently, the drawback has to be resolved by inventing a transformation that is capable of performing interpolation with a different dimension. Unfortunately, we do not have an appropriate method for solving this problem. For this reason, a fixed-order filter will serve in this study.

4.2.1 Spectral Feature

In order to meet the spectral requirements of the encoder, any source model for the LP synthesizer should have a flat spectrum. Our method for achieving the spectral feature of the glottal excitation is inspired by the appearance of the integrated excitation, in which the pulse energy around the glottal closure concentrates most of the high-frequency energy. Our formulation is given as follows:

First, we modify the third sample, $x(3)$, of the modified polynomial waveform of the integrated excitation so as to ensure the existence of a sharp peak:

$$x(3) = (1 + \max\{g(n) | n = 1, 2, \dots, [Q/2]\}^4 - 9) - 1 \quad (4-25)$$

where $g(n)$ represents the modified integrated excitation with the pitch period of 1 sample. Next, the fourth and sixth samples are changed to be

$$x(4) = \frac{x(3) + g(5)}{2} \quad (4-26)$$

and

$$x(6) = \frac{x(3) + g(5)}{2}, \quad (4-27)$$

keeping the new value of $x(3)$ so that energy at the middle frequency is enhanced. The excitation pulse is obtained by taking the difference of $g(n)$. Finally, a first-order inverse filter is applied to remove the spectral tilt of the excitation pulse. The residual excitation is defined as the glottal impulse, which will be frequently used in the rest of this dissertation.

4.2.4 Effect of Vocal Noise

Vocal noise is important for synchronizing breathy and female voices (Klatt, 1987; Fuchs et al., 1999; Klatt and Klatt, 1990; Childers and Lee, 1994). In Chapter 2, we have shown that the extracted noise exhibits the following two features. First, the noise, consisting of the individual noise and modeling errors, has a flat spectrum. Second, the magnitude of such noise near the glottal closure is higher than that at the other places.

As we also pointed out, part of the vocal noise possibly resulted from the glottal misalignment. In order to verify this possibility, we decide to use a modulation using a noise source that has a larger amplitude around the glottal closure. The noise is produced by modulating uniformly distributed white noise with a Gaussian window given by

$$W_g(t) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\sin(\pi t/T))^2}{2\sigma^2}\right) + \beta, \quad -|T/2| \leq t \leq |T/2| \quad (4-16)$$

where T is the pitch period, $[\cdot]^2$ denotes the floor function. By referencing the measurements in Chapter 2, we choose that β is 0.5 to approximate the amplitude modulation of the vocal noise for normal male subjects. While adding this noise to the synthetic solution, the amplitude of the vocal noise is adjusted to achieve a Signal-to-Noise Ratio (SNR) of 25dB.

For the sake of comparison, we also test another type of vocal noise with a constant modulation, of which the level is measured at the middle between two glottal closure events. Furthermore, we adopt a 50% duty cycle starting at the maximum glottal closure since it was preferred to a binary modulation (Childers and Lee, 1994). The SNR is modified to 25 dB to meet the measured level.

4.2.5 Source-Filter Interaction

Source-filter interaction has been emphasized to be important for synthesizing high-quality, natural-sounding speech (Wu et al., 1993; Allen and Ramey, 1995). In order

to develop a comprehensive source model for speech synthesis, this particular effect cannot be exploited. The interaction between the source and tract can be achieved either by using a vocal system to control the glottal impedance or by incorporating an excitation effect into a source model. Our approach falls in the latter category.

Two major effects in the glottal form, namely, *staircase* and *formant ripple*, result from the shorter-term interactions (Flanagan, 1973; Fant and Acemogluopoulou-Mia, 1982). The staircase, in general, results in a relatively slow rise. Our glottal impulse model is expected to account the staircase with adequate precision. The formant ripple, on the contrary, are expected to be degraded by this model. Due to the lack of accurate estimation, we only present methods for locating the formant ripples rather than direct modeling.

Since the ripple effect is associated with an increase in the formant bandwidth during the glottal open phase, similar results can be achieved by moving the poles of the LP filter (except for pole-zero cancellation) further spectral weighting filter in Section 4.1 & 2. The damping of an all-pole filter can be controlled by multiplying the corresponding coefficients, a_k 's, by the power of a factor α , i.e., $\hat{a}_k = \alpha a_k^2$ (Yamashita and Mochied, 1975; Tokimura et al., 1976). A value of α smaller than 1 will move the poles toward the origin and broaden the bandwidths of the poles. However, the opposite statement does not necessarily true when α is greater than 1. This is because the bandwidths are related only if the poles are moved closer to the unit circle.

One possible implementation for the ripple effect is to use two sets of LP coefficients to simulate the damping factors for two different glottal phases. We apply the normal LP coefficients during the glottal closed phase and switch to the LP coefficients with a larger damping (i.e., $\alpha < 1$) when the vocal folds are open. As illustrated in Figure 4-6(1), the damping of the formant energy in the synchousal vowel *ai* is quite evident.

In case the two bandwidths are narrower than those obtained by LP analysis (Wong, 1982), we instead decrease the damping of the close phase to compensate for the reduced bandwidths. Of course, this work can be accomplished by using a filter with $\alpha > 1$, but the

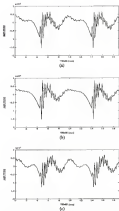


Figure 4-6: Synthetic vowel *ai* with: (a) original damping, (b) increased damping during the glottal open phase, (c) decreased damping during the glottal closed phase

Implementation of each filter requires a priori knowledge about the locations of the poles. This calls for a *peak-finding routine* which is computationally expensive and very sensitive to the quantization errors. Thus, an alternative solution that we adopt in this approach is to employ a filter, $W(z)$, to modify the excitation during the glottal-closed phase so that the synthesized speech has a similar ripple effect. The filter is given by

$$W(z) = \frac{(z - \alpha z^{-1})(1 - \sum_{k=1}^L a_k z^{-k})}{(z - \beta z^{-1})(1 - \sum_{k=1}^L a_k \beta^k z^{-k})} \quad (4-25)$$

where the values for α and β are 0.8 and 0.7, respectively.

Although $W(z)$ takes effect only in the glottal-closed phase, the filter memory carried over from the previous frame has to be taken into account. To reduce the computation of filter memory, we impose a strict constraint that the glottal impulses for any two consecutive periods are the same. Thus, filter memory can always be derived from the present glottal impulse instead of referring to the filtering results of the previous frame. By taking advantage that the duration of the glottal closed phase is usually less than one-half of the pitch period for most voices and that the memory depth is limited by the filter order, we can perform the memory recovery process together with the ripple effect by filtering a circularly shifted glottal impulse with the glottal closure as the center (see Figure 4-75). In this manner, the filter requires no memory during the glottal-open phase and distributes the memory influence to the excitation approximately during the glottal-closed phase.

To smooth the above process, we apply a *hanning window* to the circularly shifted glottal impulse. The windowed component is fed into the filter $W(z)$ to yield the intended damping. On the other hand, the remaining part (obtained by subtracting the windowed component from the excitation pulse) is kept unchanged throughout the course of filtering operation. Because $W(z)$ is an IIR filter, we append zeros to the windowed component to form a length of one and half pitch periods such that the filter is able to release its energy

sufficiently. The released energy is absorbed by adding the filtered component to the additional half period to that in the first half period. The glottal impulse with the source-tract interaction feature is formed by accumulating the filtered and synthesized components together. Figure 4-4(b) shows the synthesized signals obtained by using such a glottal impulse. In this figure, the damping of the processed speech signal is obviously reduced in the glottal closed phase.

4.2.4 Generation of Glottal Impulse

To generate a glottal impulse, the recommended order for the implementation of the spectral fluxes, vocal noise and source-tract interaction is given as follows:

1. generate a waveform $p(t)$ from the polynomial model;
2. modify $p(t)$ and $p'(t)$ by using Eqs. (4-23), (4-24) and (4-27), respectively;
3. add white noise;
4. determine the waveform to yield an excitation pulse;
5. normalize the source-tract interaction;
6. remove the spectral tilt of the excitation pulse by spectral filtering.

Figure 4-7 illustrates the procedures above.

5.1 Glott Determination

In a speech production model the glott is a function to control the pressure variation along an utterance. Although the glott is an important factor affecting synthesis quality, unfortunately it has not received much attention from researchers. In this section we discuss methods for calculating the glott and their influences on synthetic speech.

5.1.1 Glott of Vocal Excitation A_g

For a source-filter speech production model (Fletcher, 1963), the filter output can be decomposed into two components: one results from the excitation $A_g u(t)$, and the other

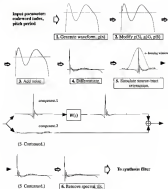


Figure 4-7. Procedure of generating a glacial impulse

made from the filter memory, $q(n)$. According to such a structure, a superposition method is often adopted for speech synthesis as to the vocal tract discretization as the foundation of pitch periods (Vahedi and Nikou, 1996). That is, for each pitch period there are two synthesis filters employed, the one holding the previous LP coefficients is in charge of the memory contribution, and the other possessing the new LP-coefficients is responsible for the current excitation.

Suppose we know that the power of the filter output has to equal that of the original signal. Given a speech segment $x(n)$ of M samples with power P_r ,

$$P_r = \frac{1}{M} \sum_{n=0}^M x^2(n) \quad (4-30)$$

And and Hansen (1971) derived the gain A_0 , by solving the following equation directly

$$P_r = \frac{1}{M} \sum_{n=0}^M [q(n) + A_0 p(n)]^2 \quad (4-31)$$

In case A_0 was negative or complex, they set A_0 as zero. The reason for choosing such a zero setting is because the power contributed by the filter memory is too small. It appears that the zero setting is just a strategy to let the filter memory die out so that the gain can increase its function. Although this zero setting seems the only solution in order to make the synthesis implementable, it definitely degrades the pitch harmonics of the synthesized speech.

Takamori et al. (1977) suggested that the memory contribution was negligible when the filter-response was sufficiently damped. Then, after increasing the damping factor of the synthesis filter, they computed the gain without considering the filter memory, i.e.,

$$P_r = \frac{1}{M} \sum_{n=0}^M [A_0 p(n)]^2 \quad (4-32)$$

In consequence, the elimination of the zero setting is at the price of possible errors due to the ignorance of filter memory.

Mallikoff (1970), on the other hand, computed the gain on the level of the driving function. By assuming that the excitation was either an impulse train for voiced speech or white noise for unvoiced speech, both of which have unity power, he computed the gain by estimating the power of the residual:

$$A_g = \left[E(E) - \sum_{p=1}^P a_p E(E) \right]^{-1/2} \quad (3-11)$$

where $E(E)$ is the autocorrelation function of the analyzed speech signal, and a_p 's are the LP coefficients. Because the gain is a by-product of the LP analysis, this method seems very elegant and straightforward. However, it leads to the following problem. Unlike the impulse, the residual signal does not have a absolutely flat spectrum. Small mismatches between the impulse and the model at low frequencies may be amplified after imposing the synthesis filter. Oftentimes, the residual errors are manifested as energy fluctuations in the synthesized speech, and a wobble-like quality will be perceived.

In CELF and MPLF coders, the part attributed to the filter memory is first removed from the analyzed speech signal (Toussaint and Ruel, 1998; Ruel and Barwell, 1999). The gain is determined by the cross correlation between the spectrally weighted speech signal and the weighted filter response of a given excitation function. This approach can automatically compensate for the residual error when the synthesized speech signal doesn't match the original very well. Unfortunately, such an approach does not suit our model because the best fit of the global impulse may still result in a large discrepancy between the original and synthesized speech signals.

From the above discussions, we see that the gain is used to replace the amplitude of the filter response of a given excitation. The spectral power P_s contains the information pertaining to the amplitude of the speech signal. To avoid the above-mentioned drawbacks, we discovered a method below to retrieve the gain from P_s .

It is noted that the speech waveforms in many adjacent pitch periods are very similar, suggesting that the initial and final filter memory are nearly equal in most cases. Thus, the filter memory can be approximated via the filtering operation with the use of the same recursion. If the number of pitch period, m , of the underlying excitation be large enough (say $m \gg 5$), the filter memory contributed before the first period are negligible. Therefore, the zero-startup filter response in the last period can be regarded as a complete filter output. Two examples of such a filtering operation are illustrated in Figure 4-3. The gain A_T for the excitation pulse is then calculated by

$$A_T = \sqrt{\frac{P_0}{1 + \sum_{n=1}^L q(n)}} \quad (4-14)$$

where $q(l)$ is the residual filter output within the m th pitch period, and L is the length of the pitch period.

The derivation described above needs a large amount of filtering operations. An algorithm presented in the following is provided to alleviate the computational burden. Notice that the filter memory from the past frame is always accessible during the speech synthesis. We can simulate the foregoing filtering process by referring to the actual filter memory. Suppose the filter is implemented using a direct-form structure. The filter memory is, therefore, represented by the coding samples of the previous frame. As was more obviously in Figure 4-3(b), it is the division of the filter memory that removes the similarity across all the periods. Based upon this observation, we separate the filter memory, $x(k)$, into two parts: the mean value, $x_m(k)$, and the deviation, $x_d(k)$:

$$x_m(k) = \frac{1}{p} \sum_{n=0}^{p-1} x(k-n) \quad \text{for } k = 1, 2, \dots, p, \quad (4-15)$$

$$x_d(k) = x(k) - x_m(k) \quad \text{for } k = 1, 2, \dots, p. \quad (4-16)$$



(a)



(b)

Figure 4-8 Zero-current responses for two variations of five-converter cycles

We simulate the filtering process by means of an iterative procedure, which is governed by a constant variance of $\delta_0^2(\delta)$. That is, at each iteration we calibrate the amplitude of the filter memory according to the mean value and deviation of the previous results by

$$\begin{aligned} x^{(i)}(x) &= u(x) + \lambda \left[\sum_{j=1}^L x_j^{(i)} U(x + \mu) - U(x) \right] + \mu \left[\sum_{j=1}^L x_j x_j^{(i)} - \delta_0(x - \bar{x}) \right], \\ (\mu &= 1, 2, \dots, J) \end{aligned} \quad (4-77)$$

where $u(x)$ is the filter response of the current iteration, $U(x)$ is the step function, and the superscript i denotes the iteration number. During each iteration, the scaling factors λ and μ are updated by

$$\lambda = \frac{1}{J} \sum_{j=1}^J x^{(j)} - \bar{x}, \quad (4-78)$$

$$\mu = \left(\frac{\sum_{j=1}^J (x_j^{(j)} - \bar{x})^2}{\sum_{j=1}^J x_j x_j^{(j)} - \bar{x}^2} \right)^{1/2}. \quad (4-79)$$

We start with the iteration from the zero-output response using the selected global response. After proceeding the above procedure several times, the resultant signal will approach to the one derived by filtering nonrecursive global signals. The gain δ_0 is calculated as

$$\delta_0 = \sqrt{\frac{P_1}{1 + \sum_{j=1}^J U^2(x_j)^2}} \quad (4-80)$$

The observed δ_0 's sometimes exhibit large variations among adjacent pixels, which may cause perceivable energy fluctuation for synthetic speech. Unfortunately, a smoothing procedure based on δ_0 's cannot remedy such a defect because it does not alter the filter gain rate across. We solve the problem by introducing a new variable, ϕ , which is

defined as the square root of the proportion of the power remaining from the current iteration

$$g = \sqrt{\frac{\sum_{n=1}^N |x(n)|^2}{\sum_{n=1}^N |x(n)|^2 + \sum_{n=1}^N |y(n)|^2}} \quad (9-41)$$

Thus it is reasonable for us to argue that the obtained g 's vary slowly during voiced speech. We therefore apply a first order IIR lowpass filter, $0.3/\delta z + 0.7z^{-1}$, to g 's to prevent occasional large power variations across pitch periods. The final result of the gain A_2 is determined by

$$A_2 = \sqrt{\frac{P_2}{1 + \sum_{n=1}^N |x(n)|^2}} \quad (9-42)$$

Eventually, with the use of the proposed algorithm, the last cross-correlation maximization by the filtering operation is replaced by multiplications and additions. More important, this algorithm prevents the drawbacks occurred in other methods.

4.3.2 Gain of Unvoiced Excitation A_0

As mentioned in Section 4.1.4.2, we have adopted the CELF algorithm to reconstruct the unvoiced excitation. The gain A_0 is the scale factor γ corresponding to the optimum codeword. While the optimal γ provides the minimum error, it also lowers the power intensity of synthetic speech. Hence, as recommended by Zeng and Kuhn (1985), the gain A_0 is better replaced by a power match between the input signal $x(n)$ (input) with the memory excitation (reconstructed) and the filtered response of the synthetic excitation $\hat{y}(n)$, i.e.,

$$A_m = \left[\frac{\sum_{j=1}^m x_j^2(0)}{\sum_{j=1}^m y_j^2(0)} \right]^{1/2} \quad (4-43)$$

where m is the length of the subframe. In our case, $m = 50$.

4.1.1. Voicing Transition

Leading to many LP synthesizers that use a multi-mode excitation source, an inappropriate voicing decision may lead to the deterioration of synthesis quality. The problem becomes serious during the voiced/unvoiced transition since the pitch estimation is prone to error at this region. The deficiency related to the pitch estimation can be alleviated by applying a median filter or some correction method to the pitch contour or that pitch halving, doubling as well as other deviating results can be avoided. To overcome the problem of a exact voicing decision, we propose a method to smooth the voicing transition as follows:

Consider a voiced segment, $s_i(0)$, that is near to an unvoiced segment. If the voiced segment is ahead of an unvoiced segment, then we can gradually change the voicing model by

$$x_i(0) = \frac{N_f + 1 - i}{N_f + 1} s_i(0) + \frac{i}{N_f + 1} s_{unv}(0), \quad i = 1, 2, 3, \dots, N_f \quad (4-44)$$

where $s_i(0)$ is the residual speech signal, N_f is the frame length, and $s_{unv}(0)$ is an alternative version of $s_i(0)$ synthesized by using the unvoicing algorithm. If the voiced segment is located behind an unvoiced segment, $x_i(0)$ becomes

$$x_i(0) = \frac{i}{N_f + 1} s_i(0) + \frac{N_f + 1 - i}{N_f + 1} s_{unv}(0), \quad i = 1, 2, 3, \dots, N_f \quad (4-45)$$

For simplicity, $s_{unv}(0)$ is derived at the analysis stage of our program. Nevertheless, it can

also be done at the synthesis stage where the CELP analysis-by-synthesis procedure is brought in to synthesize $s_d(t)$ produced using the global analysis.

4.4 Subjective Quality Evaluation

Like many other waveform coders, the use of the proposed source model is to extract important features that are not modeled by the LP filter. However, it is important to point out that although our synthetic speech waveform is very close to the original, we do not apply any closed-loop waveform-matching criterion, nor a spectral weighting function while synthesizing voiced speech. It appears that the subjective measure on the basis of segmental SNR is not appropriate to indicate the quality of synthetic speech. For this reason, we conducted informal listening tests to assess the performance of the proposed source model as well as the LP speech synthesizer.

In addition to the training sentences, two other sentences have been tested, namely, "What day was it today?" and "Should we chase those new boys." The speech segments included those uttered by speakers not in the training group. It was found that the quality of the synthetic speech was very close to that of the original speech. If the recorded speech was played back by loudspeakers in an A-B test, listeners found it difficult to discriminate the synthetic speech from its original counterpart. For speech segments in which pitch contours were identical, the probability was approximately one-third that the synthetic speech was preferred over the original speech.

To acquire a more critical view of the excitation model, the listening tests were also carried out using a high-quality loudspeaker (Sony, MDR-1W). It was revealed that our voiced-excitation model tended to deliver more energy around the fundamental frequency, and the inverse filter could not fully reconstruct such a tendency. As a result, the synthetic speech was judged to be slightly hoarse. However, when we increased the order of the inverse filter, the synthesized quality became crisper. Because such a crispy quality was not always preferred by listeners, we did not consider the increased order as an acceptable way for this.

In contrast to the finite addition of the inverse filter, both noise and source tract interaction were considered to be responsible for the improvement of synthetic quality. The addition of noise, in general, reduced the modal nature of synthetic speech. However, the use of a different amplitude modulation did not affect the speech quality largely. This is probably because the noise power for the modal resonant modes is not too significant difference. The quality improvement due to the incorporation of source-tract interaction was noticeable in our experiments. We assume that as the combined result of raising formant resonances and attenuating the non-formant components of the glottal impulse, which contributes the dispersive nature of the formant apples and, in turn, reflects the fact that the vocal tract is somewhat inelastic. From the view of spectral shaping, the resonant effect of formant apples is considered the same as the amplitude spectrum modification introduced by Kang and Everett (1983) under adaptive condition suggested by Chen and Cerny (1987). For this reason, a half-wave filtering operation that disperses the formant apples along of the glottal closure is also recommended. In other words, the $W(t)$ that we used to modify the excitation pulse can be a zero-phase filter.

From the listening part, it was also found that modification with respect to $g(t)$, $g'(t)$ and $g(t)g'(t)$ varied the pattern of vocal fold closure. According to our experience, a different closure pattern might lead to a change of the perceived quality. Although our empirical formula for constructing the excitation pulse was a variety of voices, at present, we do not have an appropriate theory to explain this result and to optimize the closure pattern.

Breathiness was reported in some synthetic speech of female speakers, especially for females with high fundamental frequencies of voicing. Using a visual comparison presented in Figure 4-9, we observe that the synthetic speech waveforms for female voices, in general, has a discontinuity on the the synthetic excitation, which showed more rapid rising slope at the OCF and less noisy components during the glottal open phase. This leads us to suspect that the vocal fold closure pattern that is suitable for male voices may be too strong for some female voices. Likewise, the level of vocal noise for males may not be appropriate for

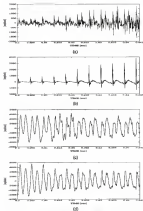


Figure 4-8 Comparison of the synthesized speech signals for a segment of voiced speech uttered by a female: (a) original speech (female), (b) synthesized speech (female), (c) original speech with added noise, (d) synthesized speech with added noise.

formants. Essentially, an ideal glottal impulse should possess certain characteristics in retaining the periodic pulses and vocal noise while preserving the segmental positions, the spectral mix of the harmonic spectrum and the continuity of the fundamental component.

Roughness was also acoustically perceived as a degradation of formant synthetic speech. Since the pitch irregularity has been considered to be an important correlate of roughness, the following two results indicate the imperfection of our GCF identification algorithm, which relied on a sharp negative peak to achieve a proper initialization and consistent similarities of adjacent pitch periods to capture the next GCFs. Because the peak was masked by nonstationary turbulence, the localized GCFs resulted in a damping effect at the later stages that even our pitch processing procedure could not fully reconstruct.

Other perceptually distortions occurred in segments containing fricatives and nasal consonants. This implies that our excitation model can only partially replicate the spectral noise (un-dominated). Because the observed glottal characteristics of the nasals are not significantly different from that of the vowels, by inference, nasal nasals are not necessarily impaired in the codetalk training. This inference has been further confirmed by testing the glottal codetalk trained without nasals. No significant degradation was found for synthesized speech using such a codetalk.

CHAPTER 5 CONCLUDING REMARKS

5.1 Summary

We confronted several problems in the first phase of this research. Attempts were made to verify the relationship between the residual signal and the glottal flow waveform. We concluded that the vocal characteristics could be retrieved from the integrated residual, which quantified the differentiated glottal flow. Also, within the source-filter theory, we proposed a comprehensive speech model that well encompasses acoustic features previously used in speech synthesis. The role of each model parameter was investigated in the context of the acoustic measures. Then, by making use of the LP analysis with the aid of EOG signals, we proposed methods for isolating and extracting the acoustic features. In particular, the perturbations of vocal sources were decomposed into low-frequency drifts and wideband noise, where the latter was extracted by using a DFT method and later applied to derive the *Singer* and *Shimmer* defined by Eqs. (3-13) and (3-14). The glottal spectral tilt was estimated using LP analysis of speech signals. While the retrieval of the spectral characteristics was performed by inverse filtering, the glottal phase was described by the envelopes of the integrated residual and a novel measure called abruptness index. The vocal noise was extracted from the integrated residual using wavelet domain approach and examined in three aspects, i.e., the relative power level, the amplitude modulation, and the noise spectrum. All the above mentioned feature extraction methods were demonstrated using sustained vowels/V's of three voice types (male, vocal fry and breathy voice) as examples. The outcomes of these acoustic measures were carefully investigated and we have reached the following conclusions:

(1) As listed in Table 3-1, the distributions of \hat{W} -factor and \hat{W} -biasness for three voice types generally agreed with other researchers' results. More important, these results substantiated our assumption that the perturbation noise-specified a Gaussian distribution in which the standard deviation was sufficient to characterize the statistical property. If we consider the quality in a broader perspective, the gross pitch and intensity variations of speech signals actually transmit linguistic messages for also non-linguistic information such as intention, emotional status, and speaker's identity, etc. In order to synthesize natural speech, an accurate and faithful replication of these variations is necessary.

(2) The turbulent noise in breathy voice was perceptually distinctive and acoustically discriminable from that in the other two voice types. This underscores the need for a vocal noise model in the source model. The noise spectra for different phonemes were fairly flat and therefore a white noise was suitable for modeling the vocal noise. On the other hand, although the amplitude modulation of the vocal noise generally resembled the magnitudes of the integrated excitation, we proved that this noise could result from the phase misalignment.

(3) The estimation of the glottal spectral tilt using LP analysis on the speech signal was tested with satisfactory results. In addition to normal inspection of the magnitude spectra of modeled filter, a simple comparison can also be carried out by using the filter coefficients of the underlying filter. The spectral tilt are moderate, relative flat, and steep for modal, vocal fry and breathy voice, respectively.

(4) The glottal phase characteristics did not show any significant relation across different voice types, suggesting no general rules for modeling the phase characteristics for different voice types. The dispersion index, in contrast, showed great potential for discriminating voice types, because the associated measures for each voice type are highly well-clustered and well-separated from one another.

The above results provide a general idea of glottal variability. More extensive investigations are needed to establish the statistical significance between model parameters

and vocal quality. The LP analysis appears to be capable of extracting the vocal source properties as well as the formant patterns. Thus, it is reasonable for us to expect that a high quality LP synthesizer is achievable if the acoustic features are accurately estimated and faithfully reproduced.

In the second phase of this research, we were interested in the design of a high-quality natural-sounding LP synthesizer. In Chapter 3, we presented an excitation model to simulate the voiced excitation by the global impulses and the unvoiced excitation by the innovation sequences. These two types of excitations were further formalized as two codebooks geared to an all-pole filter, of which the coefficients are estimated using the orthogonal covariance method. Schemes for speech analysis and synthesis were discussed in Chapter 4. Experiments with this new model and processing schemes demonstrated the feasibility of producing natural-sounding speech. In addition to source modeling, we believe efforts that lead to such accompanying works include the methods and algorithms perfecting the OCF identification, codeword matching, piece-wise LP interpolation, global pitch smoothing, spectral adjustment, source-mass interaction and gain determination. These are either involved for the first time in the literature or have had some modifications. Our achievements can be appreciated by appraising the quality of synthetic speech.

2.2 Possible Improvements

Though our LP synthesizer has been tested with fairly high success, there is still room for extension. Several possible improvements are suggested as follows.

2.2.1 Extension of Vocal Mode

Our source excitation algorithm was impeded by the difficulty of phase-mesh (pitch). While the pitch delay is always estimated as integer multiples of the sampling (or resampling) interval, a possible method for overcoming this drawback is the use of a pitch predictor, but it not only provides the necessary interpolation but also maintains the coherence between

the analyzed signals. In general, the number of filter taps need not be too many and the estimated coefficients can be easily obtained by measuring the mean squared error between the two signals. However, since the noise must be measured at the level of source variation, more studies should be made concerning the effect of the segmental order of the pitch predictor and the current filter. Furthermore, as we already pointed out in Section 2.3, there are two types of noise presented in the random signal, namely, the noise associated with the speech variation and that with the surface irregularities. One may consider how to decompose the random signal into two such components, identifying the pitch predictor to measure the speech or surface variations separately. This will allow us to enlarge our view of the vocal noise.

3.1.1 GCI Identification

The improvement of performance and reliability of the GCI identification algorithm becomes an urgent requirement for high quality speech synthesis. In this article, we located the GCI's by first choosing the largest negative peak of the integrated residue as a frame as a reference mark and then searching for the other peaks by a maximum correlation approach. The residual synthetic speech follows probability distribution emerging from the inaccurate pitch identification. Thus, we have to rely on some conditional procedure to rectify some error pitch transitions. It was reported that the GCI identification could achieve good performance if the maximum correlation approach was directly applied to the speech signal (Cheng and O'Grady, 1989). Although our experiments with Cheng's approach did show some promising results, this approach has to be further refined before it can function successfully.

3.1.2 Formant Source

Meanwhile, we are concerned with the excitation function for normal voiced, lowering weak intensity for the modification of the glottal waveform. For the utterance with elongated

glottal open phase, the increased air turbulence will certainly perturb the accuracy of the linear filter. Therefore, not only is the primary excitation pulse train distorted, but also there are other spurious components. It is obvious this type of excitation cannot be properly described by a sharp glottal impulse. Also, based upon our perceptual impression, we believe that the strong excitation pulses should at least be partially responsible for the hoarse characteristics of female speechless speech. One may think about ameliorating such defects by raising the vocal noise. According to our experiment, adding noise did increase the loudness but could not actually reduce the volume. We therefore have to resort to other means. Several efforts in the past had been directed to designing excitation signals with low peak factors and flat magnitude spectra (Schroeder, 1970; Rabiner and Crandall, 1975). Apparently, a logical follow up of this research will be to consider designing a different codebook or developing processing schemes that control the peak factor and the vocal noise as well.

5.1.4 Ripple Effect

The ripple effect is known to be important for the improvement of speech quality, but it can only be measured empirically in our experiments due to the lack of an efficient method for measuring the formant damping. Without the proper amount of the formant ripple that should be incorporated into the speech synthesis we are able to produce natural sounding speech? If we believe that the damping in formant damping only occurs at the transition from open-to-closed and/or closed-to-open glottis, then the utilization of techniques developed for fast formant tracking (Ting and Chaffin, 1980) and for the estimation of exponentially damped sinusoids (Pardollinsky and Taffel, 1983) may offer answers to this question. More studies are needed to decide how to control the damping factors and how these parameters affect the quality. Even though the estimation procedures may be computationally prohibited from practical use, we would at least gain some qualitative description that might characterize different groups of speakers.

3.2.3 Sampling Resolution

Two approaches are considered useful to improve the sampling resolution. One is the use of a multiple-rate pitch predictor that we discussed in Section 3.2.1. The other is the fractional interpolation technique discussed in Section 2.3.4. We believe the interpolated values of the signal will be able to represent the voiced speech more accurately and achieve an improvement of the synthetic quality for female speakers.

3.2.4 Spectral Distortion

In the proposed LF speech synthesizer, the spectra of speech signals are represented by an all-pole model. This model, however, is not suited for female and consonants, for the spectral envelopes of such sounds exhibit dips (zeros) between peaks (poles). Though our source excitation partially compensates for the absence of spectral zeros by exploiting the adaptiveness of random searching in Autoregressive Moving Average (ARMA) model, it seems to be even more attractive and well-suited for improving the quality of synthetic speech when spectral zeros are perceptually important (Noll and Schroeder, 1978; Aktenen and Karjal, 1989).

5.3 Applications

There are at least three areas where the techniques developed in this research are applicable.

5.3.1 Quality Measure

In this research we have postulated that a comprehensive speech production model should be composed by a complete set of acoustic functions. Thus, by utilizing the feature extraction techniques developed in Chapter 2, researchers would be able to determine a perceptually objective quality (or diagnostic) measure in a full context by verifying the relationships between the psychoacoustic attributes and the model parameters. Two possible

approaches that may achieve the described objective measures are: (I) statistical analysis based on large amounts of data, and (II) the stimulus-by-synthesis procedure.

In addition to assessing the speech quality (or severity), the duration measure can be used to study the speaker-discriminating acoustic features so that the speaker identification techniques can be advanced. Such measures can also be used to study the way the human processes the acoustic signals so as to improve the speech recognition techniques. Furthermore, for any linguistic pathology that is acoustically perceptible, the duration measure can provide insight toward the classification of the pathology.

3.3.2 Speech Coding

We have not spent much effort dealing with the process of quantization and coding, but one may anticipate the superiority of the proposed model in speech coding. Some advantages can be easily figured out by comparing the size of codeword for the various methods. The size of the global codeword required in our model is far less than that used by the Deltic & Kish voice coder (Campbell et al., 1988). In addition, only one codeword is involved to characterize the vocoder input for each frame. This will greatly reduce the demand of bit allocation and accelerate the processing speed.

3.3.3 Voice Conversion

Because every codeword represents a different pattern of the global phase characteristics, the codeword searching can be considered as a process of monitoring and tracing the phase variation. This implies that the global codeword can also be used to study the phase properties. As an example, we have applied one single-codeword to synthesize an entire sentence while keeping other parameters unchanged. The synthesized speech was intelligible but not as natural as the one synthesized by using the selected codeword sequence. Eventually, the poor quality is caused by the fixed global phase. This fact also indicates another possible application of this source model which is adaptable by other waveforms.

context. That is, our model can be used to correct the global phase characteristics. In the past, the LP synthesizers were only used to perform the prosodic and spectral modifications in the voice conversion systems (Chikami et al., 1999b; Davis and Nies, 1994; Vaheri et al., 1992). The accurate synthesis using our model will render a complete version for voice conversion.

3.1.4 Text-to-Speech Synthesis

The text-to-speech system is considered as the efficiency of converting the text to a sequence of phonetic transcription before producing the acoustic output, which includes semantic, syntactic and lexical rules that manage various intermediate transformations (Klatt, 1987). The final phonetic transcription is represented by the formant parameters (or LP coefficients), duration, pitch and intensity contours, all of which correspond the parameters of a speech production model. Thus, if the LP coefficient vectors for various phonetic segments were already registered in a codebook, we can easily apply our coefficient codebook to produce speech signals with desired quality. In fact, nowadays many LP codebooks are used for speech coding and have demonstrated high performance. It seems that the incorporation of the LP codebook into our current model is not only hypothetically feasible but also is a simple way to integrate the work of speech synthesis.

REFERENCES

- Arfken, D. J. (1996). *Elementary Fluid Dynamics* (Oxford University Press, New York).
- Att, C. (1991). "A study of vowel types and acoustic modeling: Analysis by synthesis," Ph.D. Dissertation, University of Florida, Gainesville.
- Atkeson, H. (1974). "A new look at the statistical model identification," *IEEE Trans. Acoust. Control AC-22*, 716-723.
- Abbas, M., and Kuo, S. (1989). "ACMA, model based speech coding at 16kb/s," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 146-150.
- Allen, D. R., and Rong, W. J. (1983). "A model for the synthesis of natural sounding vowels," *J. Acoust. Soc. Am.* 74(1), 58-69.
- Allen, J. B. (1977). "Short time spectral analysis and synthesis and modification by discrete Fourier transform," *IEEE Trans. Acoust., Speech, Signal Process.* 25(3), 229-238.
- Allen, J. B., and Rabiner, L. R. (1977). "A unified theory of short time spectrum analysis and synthesis," *Proc. IEEE*, 65, 1201-1244.
- Alamilli, L. B., and Silva, P. M. (1984). "Variable-frequency synthesis: An improved harmonic coding scheme," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 27.2.1-27.2.4.
- Alamilli, L., and Tribollet, J. M. (1982). "Formant coding: A low bit rate, good quality speech coding technique," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 1664-1667.
- Anandapadmasinha, T. V., and Yegnanarayana, B. (1999). "Epoch extraction from linear prediction residual for identification of voiced glottis interval," *IEEE Trans. Acoust., Speech, Signal Process.* 27(4), 305-319.
- Askenfeldt, A. G., and Hammarberg, B. (1984). "Speech waveform perturbation analysis: A perceptual-neural comparison of seven measures," *J. Speech and Hear. Dis.* 20, 58-64.
- Att, B. E., and Duvall, W. (1978). "On synthesizing natural sounding speech by linear prediction," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 44-47.
- Att, B. E., and Blamont, R. L. (1971). "Speech analysis and synthesis by linear prediction of the speech wave," *J. Acoust. Soc. Am.* 50(2), 407-409.
- Att, B. E., and Bernick, J. B. (1982). "A new model of LPC extraction for producing natural sounding speech at low bit rates," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 614-617.

- Atal, B. S., and Schrammer, M. B. (1975). "Linear prediction analysis of speech based on a pole-zero representation," *J. Acoust. Soc. Am.* 58(5), 1310-1318.
- Atal, B. S., and Schrammer, M. B. (1976). "Predictive coding of speech signals and subjective error criteria," *IEEE Trans. Acoust., Speech, Signal Process.* 24(3), 247-254.
- Berg, J. W. van den. (1988). "Myoelectric neuromagnetic theory of voice production," *J. Speech and Hear. Res.*, 1, 217-244.
- Berg, J. W. van den, Zonneba, J. T., and Doornik, P. Jr. (1977). "Clarity as resonance and the formants effect of the human larynx," *J. Acoust. Soc. Am.* 29, 626-631.
- Bergman, A., and Hsieh, P. (1989). "Cochlebank driven global pulse analysis," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 33-38.
- Berman, M. G. (1979). "Estimation of global volume velocity by the linear prediction inverse filter," Ph.D. Dissertation, University of Florida, Gainesville.
- Bladen, B. A. W., and Lutfiyya, B. (1981). "Modeling the judgments of voiced quality differences," *J. Acoust. Soc. Am.* 69(3), 1414-1422.
- Burn, A., Gray, A. H. Jr., Gray, R. M., and Markel, J. D. (1988). "Speech coding based upon vector quantization," *IEEE Trans. Acoust., Speech, Signal Process.* 36, 362-374.
- Campbell, J. P., and Thomas, E. T. (1984). "Waveform-based classification of speech with applications to the U.S. government LPC-10c algorithm," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 403-408.
- Campbell, J. P., Welch, V. C., and Thomas, T. E. (1985). "An expandable open-protected 4800-bps CELP codec," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 716-718.
- Casper, B., and Atal, B. S. (1977). "Role of multi-pulse excitation in synthesis of natural-sounding voiced speech," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 2388-2394.
- Chenbin, S., and Lin, W. C. (1974). "Experimental comparison between stationary and nonstationary formulations of linear predictors applied to voiced speech analysis," *IEEE Trans. Acoust., Speech, Signal Process.* 22(4), 400-413.
- Chen, J. H., and Gerstoft, A. (1977). "Band-pass vector LPC speech coding at 4800bps with adaptive coding," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 2163-2168.
- Cheng, Y. M., and G. Stangorou, D. (1985). "Assessment and reliable estimation of global cluster instant and period," *IEEE Trans. Acoust., Speech, Signal Processing.* 33(12), 1808-1813.
- Childen, D. G., and Fox, K. S. (1992). "Detection of laryngeal function using speech and cinematographic data," *IEEE Trans. on Biomedical Eng.* 39(1), 19-25.
- Childen, D. G., Hahn, M., and Lane, J. N. (1989a). "Silent and voiced/unvoiced/semi-vocative (five-way) classification of speech," *IEEE Trans. Acoust., Speech, Signal Process.* 37(11), 1731-1734.

- Childers, D. G., and Lang, J. N. (1984). "Electroglottography for laryngeal function measurement and speech analysis," *IEEE Trans. on Acoustical Eng.* 31(12), 897-917.
- Childers, D. G., and Lee, C. K. (1991). "Vocal quality factors: Analysis, synthesis, and perception," *J. Acoust. Soc. Am.* 90, 2394-2419.
- Childers, D. G., and Wu, K. (1990). "Quality of speech produced by analysis-synthesis," *Speech Commun.* 9, 97-117.
- Childers, D. G., Wu, K., Marks, D. M., and Yegnanarayana, B. (1989b). "Voice conversion," *Speech Commun.* 8, 143-158.
- Childers, D. G., Yeo, J. J., and Krishnamurthy, A. (1991). "Spectral analysis: A6, M6, A6/M6," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. Workshop on Spectral Estimation*, 2.2.1-2.2.18.
- Childers, D. G., Yegnanarayana, B., and Wu, K. (1988). "Voice conversion: Factors responsible for quality," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 341-344.
- Cohen, R. H. (1973). "Vocal intensity in the modal and falsetto registers," *Folia Phonetica* 23, 81-99.
- Emmberg, M. D., and Wang, D. Y. (1978). "Development of a 4.5-P 4 kHz RLLP vocoder," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 554-557.
- Endsley, J. R. (1982). "Evaluation of laryngeal dysfunction based on features of an acoustic estimate of the glottal waveform," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 159-162.
- Endsley, J. R., and Anderson, D. J. (1983). "Automatic classification of laryngeal dysfunction using the roots of the digital inverse filter," *IEEE Trans. on Acoustical Eng.* 31, 114-121.
- Dudley, H. (1939). "The vocoder," *Bell Labs Res.* 18, 123-136.
- Falkman, L., Childers, D. G., and Marks, D. M. (1990). "Acoustic correlates of vocal quality," *J. Speech and Hear. Res.* 23, 289-306.
- Fatt, G. (1969). "The acoustics of speech," *Proc. Int. Conf. on Acoust.* 148-161.
- Fatt, G. (1968). *Acoustic Theory of Speech Production* (Mouton, Paris).
- Fatt, G., and Ananthapadmanabha, T. V. (1942). "Transmission and superposition," *Speech Trans. Lab.-Q (Prog. Status Rep. (Royal Institute of Technology, Stockholm, Sweden)* 2-3, 1-17.
- Fatt, G., Lipmanowicz, J., and Lau, Q. (1983). "A four-parameter model of glottal flow," *Speech Trans. Lab.-Q (Prog. Status Rep.* 4, 1-11.
- Fatt, G., and Lau, Q. (1988). "Frequency domain interpretation and derivation of glottal flow parameters," *Speech Trans. Lab.-Q (Prog. Status Rep.* 3-3, 4-21.
- Flemming, J. L. (1973a). *Speech Analysis, Synthesis and Perception* (Springer-Verlag, New York), 2nd ed.

- Flanagan, J. L. (1972b). "Voices of man and machines," *J. Acoust. Soc. Am.* 51(2), 1575-1587.
- Flanagan, J. L., and Golden, R. M. (1966). "Phase vocoder," *Bell Syst. Tech. J.*, 45, 1493-1580.
- Papoulis, L., and Pongyoyot, M. (1986). "Proposal and evaluation of models for the global source mechanism," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 1609-1610.
- Pikarevic, T., El-Amoufy, A., and Horpa, J. (1989). "A new index for evaluation of the turbulent noise in pathological voice," *J. Acoust. Soc. Am.* 85, 1189-1195.
- Port, N. (1985). *Digital Speech Processing, Synthesis, and Recognition* (Macmillan, New York).
- Quatieri, C. F., Mowen, J. R., and Kamei, M. M. (1992). "Adaptive vocal tract production coding," *IEEE Trans. Signal Process.* 40(5), 1317-1328.
- Gold, C. (1988). "Voice source dynamics in connected speech," *Speech Trans. Lab. -Q, Prog. Status Rep.* 3-3, 173-189.
- Gold, C. (1989). "A preliminary study of acoustic voice quality," *Speech Trans. Lab. -Q, Prog. Status Rep.* 4, 9-21.
- Gree, R. M. (1964). "Voice-quality-index," *IEEE ASSP Magazine* (April), 4-29.
- Griffin, D. W., and Lee, J. S. (1988). "Multiband excitation vocoder," *IEEE Trans. Acoust., Speech, Signal Process.* 36(5), 1223-1235.
- Gulow, B., Morys, M., and Card, B. (1976). "A voice source taking account of coupling with the supraglottal cavity," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 47-50.
- Haugen, J., Nilsen, H., and Hansen S. O. (1982). "Improvements in 1.4-ks high-quality speech coding," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* B-145-B-148.
- Hochberg, P. (1981). "Acoustic-coupled source-matched vocoder," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 285-288.
- Hochberg, P. (1986). "High quality global LPC vocoding," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 685-688.
- Hochberg, P. (1988). "Phase-compensation in all-pole speech analysis," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 109-112.
- Kamathakr, H., Hansen, F. A., and Walker, H. (1985). "Perceptually-based linear predictive analysis speech," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 306-312.
- Kita, S., Iwamoto, S., Hirose, M., Mizushima, H., and Kohno, Y. (1978). "Acoustical analysis for voice disorders," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 613-616.
- Politzer, R. E., and Weinberg, B. (1961). "Estimation of volume velocity waveform properties, A review and study of some methodological assumptions," in *Speech and Language: Advances in Basic Research and Practice*, edited by N. Luce (Academic P., New York), pp. 411-423.

- Hirooka, M., Kikuchi, Y., Ueda, H., Tanaka, S., and Tanabe, M. (1994). "Acoustic-intensity analysis of normal and hoarse voices," *J. Acoust. Soc. Am.* 95, 1641-1651.
- Holmes, R. (1994). "The vocal register," *J. Phys.* 2, 125-144.
- Holmes, J. N. (1975). "The influence of glottal waveform on the intensity of speech from a parallel formant synthesizer," *IEEE Trans. Audio Electroacoust.* AU-23(3), 294-303.
- Holmes, J. N. (1982). "Formant synthesizers, cascade or parallel," *Speech Commun.* 2, 251-274.
- Jerkis, R. R., Antoniano-Daniels, N., and Macdonald, L. (1987). "Digital source filtering for linguistic research," *J. Speech and Hear. Res.* 30, 127-129.
- Kahn, M., and Gans, R. (1982). "The effects of five voice characteristics on LPC quality," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 310-343.
- Kang, Q. S., and Stevens, K. (1982). "Improvement of the excitation source in the narrow-band linear prediction vocoder," *IEEE Trans. Acoust., Speech, Signal Process.* 30(3), 337-346.
- Katano, S. (1988). "Glottal waveform parameters for different speaker types," *Speech Trans. Lab. -Q. Prog. Status Rep.* 3-4, 41-47.
- Kawaya, H., Ogawa, S., and Kohuchi, Y. (1984a). "An acoustic analysis of pathological voice and its application to the evaluation of laryngeal pathology," *Speech Commun.* 3, 171-181.
- Kawaya, H., Ogawa, S., Matsumi, K., and Kohuchi, Y. (1984b). "Normalized noise energy as an acoustic measure to evaluate pathological voice," *J. Acoust. Soc. Am.* 84(3), 1329-1334.
- Kawaya, K. (1981). "Quantitative evaluation of the noise level in the pathological voice," *Folia Phoniatrica*, 33, 103-124.
- Klan, D. R. (1988). "Software for a cascade/parallel formant synthesizer," *J. Acoust. Soc. Am.* 83(3), 991-995.
- Klan, D. R. (1987). "Various source-to-speech conversions for English," *J. Acoust. Soc. Am.* 82(3), 757-769.
- Klan, D. R., and Klan, L. C. (1990). "Analysis, synthesis, and perception of voice quality variations among female and male talkers," *J. Acoust. Soc. Am.* 87(2), 830-857.
- Kleijn, W. B., Krimholtz, D. J., and Rasmussen, R. H. (1990). "Fast methods for the CELP speech coding algorithm," *IEEE Trans. Acoust., Speech, Signal Process.* 38(6), 1330-1342.
- Kolka, Y. (1985). "Vowel sympathetic modulation in patients with laryngeal diseases," *J. Acoust. Soc. Am.* 85(4), 1339-1346.
- Kudo, T., and Mikiishi, J. (1975). "Application of source filtering for detecting laryngeal pathology," *Annals of Otology, Rhinology and Laryngology* 84(1), 117-124.

- Kanamori, K. (1981) "Threshold bounds in SVQ and a new recursive algorithm for order selection in AR models," *IEEE Trans. Signal Process.* 29(5), 1119-1127.
- Konikominas, K., and Tan, K. (1990) "Statistical analysis of effective singular values in matrix rank determination," *IEEE Trans. Acoust., Speech, Signal Process.* 38(1), 317-331.
- Krishnaswamy, A. R., and Chidambaram, G. C. (1985) "Two-channel speech analysis," *IEEE Trans. Acoust., Speech, Signal Process.* 33(4), 130-143.
- Kruis, P., and Aal, B. S. (1990) "Pitch prediction with high temporal resolution," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 661-664.
- Kysenbary, H. (1986) "A pitch-synchronous analysis/synthesis system to independently modify formant frequencies and bandwidths for voiced speech," *Speech Commun.* 3, 311-326.
- Kwon, S. Y., and Goldberg, A. J. (1984) "An enhanced LPC vocoder with no overflows/underflows," *IEEE Trans. Acoust., Speech, Signal Process.* 32(4), 455-458.
- Labrent, A. L., and Chidambaram, B. G. (1981) "Modeling vocal disorders via formant synthesis," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 585-588.
- Lewis, F., and Hanson, R. (1981) "Describing the normal voice," in *Evaluation of Speech in Psychosurgery* ed by J. Early (Grun and Stroun, New York), pp. 51-79.
- Lehman, P. (1961) "Pitchbent in vocal pitch," *J. Acoust. Soc. Am.* 33(5), 587-603.
- Lehman, R., and Blount, S. E. (1988), *Speech Physiology, Speech Perception, and Articulatory Phonetics* (Cambridge U. P., New York).
- Leeds, T., Berry, A., and Ong, B. M. (1981) "An algorithm for vector quantizer design," *IEEE Trans. Commun.* COM-29(1), 54-64.
- Ma, C. K., and Chan, C. K. (1981) "Maximum-likelihood method for image vector quantization," *Electronics Letters*, 17(26), 1172-1173.
- Markel, J. (1975) "Linear prediction: A tutorial review," *Proc. IEEE*, 63, 561-580.
- Markel, J., Yipmanthan, R., Schwartz, R., and Higgins, A. W. F. (1978) "A mixed-source model for speech compression and synthesis," *J. Acoust. Soc. Am.* 64(5), 1371-1381.
- Markel, J. D., and Ong, A. H. (1976) *Linear Prediction of Speech* (Springer-Verlag, New York).
- Masouk, M., and Soudov, Y. (1982) "A new approach to the determination of the global iteration," *IEEE Trans. Acoust., Speech, Signal Process.* 30(5), 616-622.
- McAulay, R. J., and Quatieri, T. F. (1984) "Magnitude-only reconstruction using a sinusoidal speech model," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 27.6 1-27.6.4.
- McCabe, A. V., and Rasmussen, T. F. (1981) "A new mixed excitation LPC vocoder," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 895-898.

- Mikelsänt, P. (1986). "Global inverse filtering by joint constraint of an AR system with a linear input model." *IEEE Trans. Acoust., Speech, Signal Process.* 34(1), 28-42.
- Mikelsänt, P. H. (1987). "Least mean square measures of vowel perturbations." *Journal of Speech and Hearing Res.* 30, 229-236.
- Mintes, R., and Engelenstein, M. (1977). "Study of variations in the male and female glottal wave," *J. Acoust. Soc. Am.* 62(4), 981-993.
- Moss, G. P. (1974). "Observation on laryngeal function, laryngeal behavior, and voice." *Annals of Otology, Rhinology and Laryngology* 83, 353-367.
- Moss, H., Maréchal, T., Wagaonara, K., Palade, H., Tiliapava, B., Fujitaka, T., and Kawan, S. (1987). "Analysis of hoarse voices using the LPC Method," in *Laryngeal Function in Phonetics and Experimentation*, edited by T. Basso-C. Sankin, and K. Hanna-Calgary-Hill Press, Boston, MA, pp. 463-474.
- Nag, T., and Whiting, S. (1983). "Power spectrum estimates via orthogonal transformation," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 2375-2378.
- Nyquist, A., and Tisserand-Bouazery, A. (1982). "Maximum-entropy estimation technique for image coding vector quantizer design," *Electronics Letters*, 24(3), 275-276.
- Oppenheim, A. V. (1989). "A speech analysis-synthesis system based on homomorphic filtering," *J. Acoust. Soc. Am.* 85(2), 458-465.
- Oppenheim, A. V., and Wilking, A. S. (1983). *Signal and Systems* (Prentice-Hall, Englewood Cliffs, NJ).
- Pao, Y.-H. (1985). *Adaptive pattern recognition and neural networks* (Addison-Wesley, New York).
- Perkowitz, S., and Cohen, C. R. (1983). "On automatic estimation of articulatory parameters in a text-to-speech system," *Computer Speech and Language* 4, 10-15.
- Perkowitz, S., and Telle, D. W. (1987). "Stochastic systemwide modeling of natural speech," *IEEE Trans. Acoust., Speech, Signal Process.* 35(7), 1241-1249.
- Pinto, N. E., Chaffin, D. G., and Lalwani, A. L. (1989). "Formant speech synthesis: Improving production quality." *IEEE Trans. Acoust., Speech, Signal Process.* 37(12), 1831-1842.
- Pinto, N. E., and Tsao, L. E. (1990). "Influence of perturbation measures in speech synthesis." *J. Acoust. Soc. Am.* 87(3), 1278-1288.
- Piscot, D. E., and Blumstein, S. (1986). "Perceptual evaluation of MITalk: The MIT synthesized text-to-speech system," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 573-575.
- Piscot, D. E., Blumstein, H. C., Liao, P. A., and Schvach, E. C. (1983). "Perceptual evaluation of synthetic speech: Some considerations of the non-synthetic interface," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 528-534.
- Pinto, P. F. L., Slevin, E. H., and Chaffin, D. G. (1993). "Optimization of acoustic-to-articulatory mapping," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 11-10-11-16.

- Pross, P. J. (1990). "Male and female voice source characteristics: Inverse filtering results," *Speech Commun.* 3, 163-177.
- Pross, A. R., Montgomery, R. R., and Hawkins, D. R. (1987). "An evaluation of random formants as correlates of vowel identity," *J. Commun. Disorders* 20, 103-113.
- Rabiner, L. R., and Crochiere, B. E. (1979). "On the design of all-pass signals with peak amplitude constraints," *Bell Syst. Tech. J.* 58(4), 963-987.
- Rabiner, L. R., and Salamon, D. W. (1978). *Digital Processing of Speech Signals* (Prentice-Hall, Englewood Cliffs, NJ).
- Robinson, D., and Elliott, R. (1958). "A determination of the equal-loudness relation for pure tones," *Br. J. Appl. Physics* 7, 148-149.
- Ross, B. C., and Ramwell, P. F. JR. (1980). "Design and performance of an analyzer for synthesis class of parametric speech coders," *IEEE Trans. Acoust. Speech, Signal Process.* 28(5), 1409-1423.
- Rosenberg, A. E. (1971). "Effect of glottal pulse shape on the quality of natural vowels," *J. Acoust. Soc. Am.* 49(2), 583-590.
- Rosenberg, M. (1982). "Acoustic interaction between the glottal source and the vocal tract," in *Vocal Fold Physiology* edited by K. N. Stevens, and M. Hirano (Shirai Tokyo Press), pp. 305-321.
- Sathian, A. R., Rosenberg, L. R., Kallman, L. R., and McGonigal, C. A. (1978). "On reducing the form in LPC synthesis," *J. Acoust. Soc. Am.* 63(3), 919-924.
- Savva, M., and Nam, I.-R. (1990). "Vowel prosody transformation," *Digital Signal Processing* 1, 107-113.
- Schaefer, R. W., and Rabiner, L. R. (1971). "A digital signal processing approach to interpolation," *Proceedings of the IEEE*, 59, 690-702.
- Schwaninger, J. (1982). "Quantitative evaluation of the discrimination performance of acoustic features in detecting laryngeal pathology," *Speech Commun.* 1, 269-282.
- Schwaninger, J. (1989). "Noise in sustained vowels and isolated monosyllables produced by dysphonic speakers," *Speech Commun.* 3, 61-79.
- Schroeder, M. R. (1978). "Synthesis of low-pass-filter signals and linear response with low autocorrelation," *IEEE Trans. Inform. Theory* 27-18, 83-89.
- Schroeder, M. R., and Atal, B. S. (1940). "Code-excited linear prediction (CELP): High-quality speech at very low bit rates," *Proc. IEEE Int. Conf. Acoust. Speech, Signal Process.* 931-940.
- Schulthaus, M. and Luecke, A. (1989). "On the performance of CELP algorithms for low rate speech coding," *Proc. IEEE Int. Conf. Acoust. Speech, Signal Process.* 152-155.
- Schwarzer, R. T., and Chao, C. D. (1985). "Characterization of aperiodicity in steady periodic signals," *Proc. IEEE Int. Conf. Acoust. Speech, Signal Process.* 1161-1164.

- Scieszka, G. (1978). "Extension of the dimension of the model." *Ann. Stat.* 6, 461-464.
- Singhal, S., and Atal, B. S. (1977). "Amplitude optimization and pitch prediction in multipulse coders." *IEEE Trans. Acoust., Speech, Signal Process.* 25(3), 317-323.
- Solich, A. M., and Childers, D. G. (1983). "Laryngeal excitation using distortions from speech and the electrophonograph." *IEEE Trans. on Biomedical Eng.* 30(11), 755-759.
- Sorensen, D., and Houti, Y. (1984). "Dimensional perturbation factors for jitter and distortion." *J. Commun. Disorders* 17, 143-154.
- Srinivasan, T. V. (1984). "Modelling LPC-analysis by components for good quality speech coding." *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 171-174.
- Solomon, R. A., LeCuney, J. L., and Parnis, J. W. (1984). "Decomposition of the LPC excitation using the three-term function." *IEEE Trans. Acoust., Speech, Signal Process.* 32(5), 1529-1541.
- Teager, S., and Houtman, T. (1988). "A global waveform model for high-quality speech synthesis." *J. Acoust. Soc. Am.* 83, 8112.
- Ting, Y. L., and Childers, D. G. (1986). "Speech analysis using the weighted-least-square least squares algorithm with a variable forgetting factor." *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 389-392.
- Titus, I. R., Myers, W., and Scherer, R. C. (1987). "Some technical considerations in voice perturbation measurements." *J. Speech and Hear. Res.* 30, 203-208.
- Tokumoto, Y., Ishikawa, F., and Hashimoto, S. (1979). "Spectral recording technique in PARCOR speech analysis-synthesis." *IEEE Trans. Acoust., Speech, Signal Process.* 28(5), 587-594.
- Tou, J. T. (1978). "DYWIDC — A dynamic optimal cluster-seeking technique," *Int. J. Comput. Inf. Sci.* 16(3), 341-347.
- Tou, J. T., and Gonzalez, R. C. (1974). *Pattern Recognition Principles* (Addison-Wesley, New York).
- Toussaint, L. M., and Atal, B. S. (1986). "Efficient search procedure for selecting the optimum excitation in stochastic coders." *IEEE Trans. Acoust., Speech, Signal Process.* 34(3), 378-386.
- Toussaint, L. M., Mungara, J. S., and Roberts, C. M. (1983). "CELP and sinusoidal coders: Two solutions for speech coding at 4.8-9.6 kbps." *Speech Commun.* 3, 389-400.
- Un, C. K., and Magill, D. T. (1975). "The residual-excited linear prediction vocoder with excitation rate below 9.6 kbps." *IEEE Trans. Commun. COM-23(11)*, 1468-1474.
- Valter, H., Hashemi, S., and Taheri, J. P. (1991). "Voice reconstruction using PCELA technique." *Speech Commun.* 11, 173-182.
- Verhelde, W., and Nijma, P. (1985). "A modified superposition speech synthesizer and its application." *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 2807-2810.

- Wangenheim, E., and Mathias, J. (1977). "Quantitative properties of summation parameters in linear predictive systems," *IEEE Trans. Acoust., Speech, Signal Process.* 25(2), 309-311.
- Wang, S., Sakag, A., and Gersho, A. (1984). "Auditory description measures for speech coding," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 483-486.
- Wax, M. (1988). "Order selection for AR models by predictive linear systems," *IEEE Trans. Acoust., Speech, Signal Process.* 36(4), 543-548.
- Wentzold, R. W. (1983). "Laryngeal analog synthesis of harsh voice quality," *John Phonetics* 15, 341-350.
- Wells, V. L., and Bernfield, T. M. (1987). "Prediction of vocal intensity scales and source-voice types," *J. Speech and Hear. Res.* 30, 379-389.
- Wong, C.-H. (1991). "The incorporation of glottal source vocal tract interaction effects to improve the naturalness of synthesized speech," Ph.D. Dissertation, University of Florida, Gainesville.
- Wong, D. Y. (1988). "On understanding the quality problems," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 725-728.
- Wong, D. Y., and Mathias, J. D. (1988). "An excitation function for LPC systems which retains the human glottal phase characteristics," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 171-174.
- Wong, D. Y., Markel, J. D., and Ong, A. H. Jr. (1979). "Least squares glottal source filtering from the acoustic speech waveform," *IEEE Trans. Acoust., Speech, Signal Process.* 27(4), 389-395.
- Wu, J. J., Krishnamoorthy, A. K., Nank, J. R., Moore, G. F., and Childers, D. G. (1983). "Glottal opening for speech analysis and synthesis," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 1332-1335.
- Yamato, E., Gould, W., and Bass, T. (1983). "Harmonics-to-noise ratio as an index of the degree of hoarseness," *J. Acoust. Soc. Am.* 71, 1549-1558.
- Yamato, E., Suzuki, Y., Okamura, H. (1984). "Harmonics-to-noise ratio and psychoacoustic measurement of the degree of hoarseness," *J. Speech and Hear. Res.* 17, 2-6.
- Zhang, X., and Chang, X. (1990). "A new excitation model for LPC vocoder at 2.4 KHz," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 1-45-1-48.
- Zhang, B. L., and Kuo, S. B. (1989). "1600 and 7200 bps hybrid vocoder multipulse coding," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* 747-750.

BIOGRAPHICAL SKETCH

Hsueh-Hua Horng was born at Gueiguan, a town in north Taiwan, Republic of China, on January 13, 1944. He graduated from National Cheng Kung University, Tainan, Taiwan, in June, 1965 with a Bachelor of Science degree in electrical engineering. After complying with the compulsory military service, he entered the Department of Electrical Engineering, University of Florida, where he received his Master of Science degree in May, 1970. Since then, he has been a graduate research assistant under the supervision of Dr. D. G. Childers at the Manif. Machine Interaction Research Center, where his primary interest is digital signal processing with application to speech analysis and synthesis. After completing the requirements for the Ph.D. degree, he intends to return to his country and participate, going forward, in research areas such as VLSI signal processing, digital image and speech processing, as well as electronic telecommunications.

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.



Donald G. Childers, Chairman
Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.



Leon W. Cook, II
Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.



Fred L. Drake
Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.



James C. Friesen
Associate Professor of Electrical Engineering

I certify that I have read this study and that in my opinion it conforms to acceptable standards of scholarly presentation and is fully adequate, in scope and quality, as a dissertation for the degree of Doctor of Philosophy.


Mark C. E. Yang
Professor of Sociology

This dissertation was submitted to the Graduate Faculty of the College of Engineering and to the Graduate School and was accepted as partial fulfillment of the requirements for the degree of Doctor of Philosophy.

May, 1993


Richard M. Phillips
Dean, College of Engineering


Nicholas M. Lockhart
Dean, Graduate School